

ESTIMATING THE INTENSITY OF A RANDOM MEASURE BY HISTOGRAM TYPE ESTIMATORS

YANNICK BARAUD AND LUCIEN BIRGÉ

ABSTRACT. The purpose of this paper is to estimate the intensity of some random measure N on a set \mathcal{X} by a piecewise constant function on a finite partition of \mathcal{X} . Given a (possibly large) family \mathcal{M} of candidate partitions, we build a piecewise constant estimator (histogram) on each of them and then use the data to select one estimator in the family. Choosing the square of a Hellinger-type distance as our loss function, we show that each estimator built on a given partition satisfies an analogue of the classical squared bias plus variance risk bound. Moreover, the selection procedure leads to a final estimator satisfying some oracle-type inequality, with, as usual, a possible loss corresponding to the complexity of the family \mathcal{M} . When this complexity is not too high, the selected estimator has a risk bounded, up to a universal constant, by the smallest risk bound obtained for the estimators in the family. For suitable choices of the family of partitions, we deduce uniform risk bounds over various classes of intensities. Our approach applies to the estimation of the intensity of an inhomogeneous Poisson process, among other counting processes, or the estimation of the mean of a random vector with nonnegative components.

1. INTRODUCTION

The aim of the present paper is to design a new model selection procedure in a statistical framework which is general enough to cope simultaneously with the following estimation problems.

Problem 1: Estimating the means of nonnegative data. The statistical problem that initially motivated this research was suggested by Sylvie Huet and corresponds to the modeling of data coming from some agricultural experiments. In such an experiment, the observations are independent nonnegative random variables N_i with mean s_i where i varies among some finite index set \mathcal{X} . In this framework, our aim is to estimate the vector $(s_i)_{i \in \mathcal{X}}$.

Problem 2: Estimating the intensity of a Poisson process. We recall that a Poisson process N on the measurable set $(\mathcal{X}, \mathcal{A})$ with finite mean measure ν is a random measure N on \mathcal{X} such that

- for any $A \in \mathcal{A}$, $N(A)$ is a Poisson random variable with parameter $\nu(A)$;
- for any family A_1, \dots, A_n of disjoint elements of \mathcal{A} , the corresponding random variables $N(A_1), \dots, N(A_n)$ are independent.

Date: April, 2006.

2000 Mathematics Subject Classification. 62G05.

Key words and phrases. Model selection - Histogram - Discrete data - Poisson process - Intensity estimation - Adaptive estimation.

We can always assume that ν is finite by suitably restricting the domain of observation of the process. When the mean measure ν is dominated by some given measure λ on \mathcal{X} then the nonnegative function $s = d\nu/d\lambda$ is called the intensity of N . A Poisson process can be represented as a point process on the set \mathcal{X} . Each point represents the time (if $\mathcal{X} = \mathbb{R}_+$) or location of some event. For example, the successive times of failures of some machine can be represented by a Poisson process on $\mathcal{X} = \mathbb{R}_+$. The intensity of the process models the behaviour of the machine in the following way: the intervals of times on which the intensity takes large values correspond to periods where failures are expected to be frequent and in the opposite, those on which the intensity is close to 0 are periods on which failures are rare. In this statistical framework, our aim is to estimate the intensity s on the basis of the observation of N .

Problem 3: Estimating a hazard rate. We consider an n sample T_1, \dots, T_n of non-negative real valued random variables with common density p (with respect to the Lebesgue measure on \mathbb{R}_+) and assume these to be (possibly) right-censored. This means that there exists i.i.d. random variables C_1, \dots, C_n such that we actually observe the pairs $X_j = (\tilde{T}_j, D_j)$ for $j = 1, \dots, n$ with $\tilde{T}_j = \min\{T_j, C_j\}$ and $D_j = \mathbb{1}_{\{T_j \leq \tilde{T}_j\}}$. Such censored data are common in survival analysis. Typically, T_i corresponds to a time of failure or death which cannot be observed if it exceeds time C_i . Our aim, here, is to estimate the hazard rate s of the T_i defined for $t \geq 0$ by $s(t) = p(t)/\mathbb{P}(T_1 \geq t)$.

Problem 4: Estimating the intensity of the transition of a Markov process.

Let $\{X_t, t \geq 0\}$ be a Markov process on \mathbb{R}_+ with cadlag paths and a finite number of states. We distinguish two particular states, named 0 and 1, and assume that 0 is absorbant and that there is a positive probability to reach 1. Our aim is to provide an estimation of the intensity of the transition time $T_{1,0}$ from state 1 to 0. Typical examples arise when 0 means “death”, “failure”, An alternative example could be the situation where $T_{1,0}$ measures the age at which a drug addict makes the transition from soft drugs (state 1) to hard drugs (state 0). In this case we stop the chain at 0 making this state absorbing. For $t > 0$, we denote by X_{t-} the left-hand limit of the process X at time t and assume that for some measurable nonnegative function p , $\mathbb{P}(T_{1,0} \leq t) = \int_0^t p(u)du$. Note that p is merely the density of $T_{1,0}$ if $T_{1,0} < +\infty$ a.s. which we shall not assume. Our aim is to estimate the transition intensity s of $T_{1,0}$ which is defined for $t > 0$ by $s(t) = p(t)/\mathbb{P}(X_{t-} = 1)$.

For pedagogical reasons mainly, since it has already been extensively studied and can therefore serve as a reference, it will be interesting to consider also the much more classical

Problem 0: Density estimation. It is the problem of estimating an unknown density s from n i.i.d. observations X_1, \dots, X_n with this density.

All the problems described in the above examples amount to estimating a function s mapping \mathcal{X} to \mathbb{R}_+ . For this purpose, we choose a family \mathcal{M} of partitions of \mathcal{X} and for each $m \in \mathcal{M}$ we design a non-negative estimator \hat{s}_m of s which is constant on the elements of this partition. We shall call such an estimator an histogram-type estimator. The performance of \hat{s}_m depends on both s and m . Since s is unknown, we cannot pick the partition which leads to the best estimator. To select a partition in \mathcal{M} , we shall rather use a method solely based on our data leading to some random partition \hat{m} and define our resulting estimator as $\hat{s}_{\hat{m}}$. Our objective is to design

the selection procedure in such a way that $\hat{s}_{\hat{m}}$ performs almost as well as the best estimator among the family $\{\hat{s}_m, m \in \mathcal{M}\}$.

The purpose of this paper is to describe some general setup which allows to deal with all the five problems simultaneously, to explain the construction of our histogram-type estimators \hat{s}_m , to design a suitable selection procedure \hat{m} and to study the performance of the resulting estimator $\hat{s}_{\hat{m}}$. We shall illustrate our results by numerous examples of family of partitions and target functions s of interest. For the problems of estimating the intensity of a Poisson process or a hazard rate on the line, our method provides estimators than can cope with different families of functions simultaneously, including monotone, Hölderian, or piecewise constant with a few jumps with unknown locations and sizes. In the multivariate case, we shall also provide some special method for estimating Poisson intensities with a few spikes with unknown locations and heights.

The problem of estimating s by model selection in the first four setups described above did not receive much attention in the literature with a few noticeable exceptions. Problem 1 is generally viewed as a regression problem where the mean s_i takes the form $f(x_i)$ for some design points x_i (typically f is defined on $[0, 1]$ and $x_i = i/n$). To perform model selection, one introduces a wavelet basis and performs a shrinkage of the estimated coefficients of f with respect to this basis. This amounts to selecting which coefficients will be kept. To this form of selection pertain the papers by Antoniadis, Besbeas and Sapatinas (2001), Antoniadis and Sapatinas (2001). Closer to our approach is Kolaczyk and Nowak (2004) based on penalized maximum likelihood. Unlike ours, their approach requires that the means s_i be uniformly bounded from above and below by known positive constants. For Problem 2, a similar approach based on wavelet shrinkage is developed in Kolaczyk (1999), but the reference result is Reynaud-Bouret (2003). Problems 3 and 4 amount to estimating Aalen's multiplicative intensity s of some counting process with a bounded number of jumps. The problem of non-parametric estimation of Aalen's multiplicative intensities has been considered by Antoniadis (1989) who uses penalized maximum likelihood estimation with a roughness penalty and gets uniform rates of convergence over Sobolev balls. Van de Geer (1995) considers the Hellinger loss and establishes uniform estimation rates for the maximum likelihood estimator over classes of intensities with controlled bracketing entropy. Grégoire and Nembé (2000) extend the results of Barron and Cover (1991) about density estimation to that of intensities. Wu and Wells (2003) and Patil and Wood (2004) derive asymptotic results for thresholding estimators based on wavelet expansions. All these results, apart from those of van de Geer, are of an asymptotic nature. Reynaud-Bouret (2002) introduces a model selection procedure to estimate the intensity. A common feature of these papers lies in the use of martingales techniques (apart from Grégoire and Nembé, 2000). Unlike theirs, our approach does not require any martingale argument at all.

In Section 2, we present a general statistical framework which allows to handle simultaneously all the examples we have mentioned. We also make a review of some special classes of target functions and the various families of models (partitions) to be used in our estimation procedure. The treatment of our five estimation problems is provided in Sections 4 and 5. The results presented there derive from a unifying theorem to be found in Section 6. The remainder of the paper is devoted to the most technical proofs.

In the sequel, we shall make a systematic use of the following notations: constants will be denoted by C, C', c, \dots and may change from line to line; we denote by \mathbb{N}^* the set of positive integers and we write $x \wedge y$ for $\min\{x, y\}$, $x \vee y$ for $\max\{x, y\}$ and $|m|$ for the cardinality of a set m .

2. PRESENTATION OF OUR METHOD

2.1. A general statistical framework. We consider an abstract probability space $(\Omega, \mathcal{E}, \mathbb{P})$ and a measurable space $(\mathcal{X}, \mathcal{A})$ bearing a nonnegative σ -finite measure λ . In the sequel \mathbb{E} will denote the expectation with respect to \mathbb{P} . We then consider on \mathcal{X} a nonnegative bounded random process $Y = Y(x, \omega)$, i.e. a measurable function from $\mathcal{X} \times \Omega$ to \mathbb{R}^+ , and the nonnegative random measure M on \mathcal{X} given by $dM = Y d\lambda$. Besides M , we also observe a nonnegative random measure N on \mathcal{X} which satisfies

$$(1) \quad \mathbb{E}[N(A)] = \mathbb{E} \left[\int_A s dM \right] < +\infty, \quad \text{for all } A \in \mathcal{A},$$

for some deterministic nonnegative and measurable function s on \mathcal{X} . Note that this assumption implies that N is a.s. a finite measure. Our aim is to estimate s from the observations N and M . Hereafter, we shall deal with estimators that belong to the cone \mathcal{L} of nonnegative measurable functions t on $\mathcal{X} \times \Omega$ such that $\mathbb{E} \left[\int_{\mathcal{X}} t dM \right] < +\infty$. Note that s also belongs to \mathcal{L} . To measure the risks of such estimators, we endow \mathcal{L} with the *quasi-distance* (since we may have $H(t, t') = 0$ with $t \neq t'$) H between two elements t and t' of \mathcal{L} by

$$H^2(t, t') = \int_{\mathcal{X}} \left(\sqrt{t} - \sqrt{t'} \right)^2 dM,$$

and set as usual, for $t \in \mathcal{L}$ and $\mathcal{F} \subset \mathcal{L}$, $H(t, \mathcal{F}) = \inf_{f \in \mathcal{F}} H(t, f)$. Given an estimator \hat{s} of s , i.e. a measurable function of N and Y with $\hat{s} \in \mathcal{L}$, we define its risk by $\mathbb{E} [H^2(\hat{s}, s)]$. In most of our applications, Y is identically equal to 1 in which case $M = \lambda$ is deterministic and if t and t' are densities with respect to M , H is merely the Hellinger distance between the corresponding probabilities. Only the cases of Problems 3 and 4 require to handle random measures M .

In order to define our estimators we assume that

$$(2) \quad \mathbb{P}[N(A) > 0 \text{ and } M(A) = 0] = 0 \quad \text{for all } A \in \mathcal{A},$$

a property which is automatically fulfilled when $M = \lambda$ is deterministic because of (1).

2.2. Histogram-type estimators. Let us now introduce the histogram-type estimators \hat{s}_m based on some finite partition m of \mathcal{X} . We consider the subset $\mathcal{J} = \{A \in \mathcal{A} \mid \mathbb{E}[M(A)] < +\infty\}$ of \mathcal{A} and define the *model* S_m as the set of (possibly random) nonnegative piecewise constant functions on \mathcal{X} :

$$S_m = \left\{ t = \sum_{I \in m \cap \mathcal{J}} t_I \mathbb{1}_I \mid t_I = t_I(\omega) \in \mathbb{R} \text{ for all } I \in m, \omega \in \Omega \right\} \cap \mathcal{L}.$$

We then define the histogram estimator \hat{s}_m as the element of S_m given (with the convention $0/0 = 0$) by

$$\hat{s}_m = \sum_{I \in m \cap \mathcal{J}} \frac{N(I)}{M(I)} \mathbb{1}_I.$$

Note that \hat{s}_m is a.s. well-defined because of (2). We shall, hereafter, call it the *histogram estimator* based on m .

Under suitable assumptions that will be satisfied for Problems 0, 1 and 2 (the case of hazard rates and Markov processes being more complicated), we shall prove for \hat{s}_m a risk bound of the form

$$(3) \quad \mathbb{E} [H^2(\hat{s}_m, s)] \leq C_0 \{ \mathbb{E} [H^2(s, S_m)] \} + C_P |m|,$$

where C_0 is a numerical constant and C_P depends on the problem we consider. For instance, $C_P = n^{-1}$ for density estimation and $C_P = 1$ for estimating the intensity of a Poisson process. We recover here the usual decomposition of the risk bounds into an approximation term which involves the distance of the parameter from the model and a complexity term proportional to the number $|m|$ of parameters that describe the model.

2.3. The selection procedure. Given the family of models $\{S_m, m \in \mathcal{M}\}$ corresponding to a finite or countable family \mathcal{M} of partitions m , we consider, in order to define our model selection procedure, the possibly enlarged family

$$\overline{\mathcal{M}} = \{m \vee m' \text{ for } m, m' \in \mathcal{M}\}; \quad m \vee m' = \{I \cap I' \mid I \in m, I' \in m', I \cap I' \neq \emptyset\},$$

so that $m \vee m'$ is again a finite partition of \mathcal{X} .

We shall systematically make the following assumption about the family \mathcal{M} .

H: *There exists some $\delta \geq 1$ such that $|m \vee m'| \leq \delta (|m| + |m'|)$ for all $(m, m') \in \mathcal{M}^2$.*

We then introduce a penalty function “pen” from \mathcal{M} to \mathbb{R}_+ to be described below and, for $m \neq m' \in \mathcal{M}$ we consider the test statistic

$$(4) \quad T_{m,m'}(N) = H^2(\hat{s}_m, \hat{s}_{m \vee m'}) - H^2(\hat{s}_{m'}, \hat{s}_{m \vee m'}) + 16[\text{pen}(m) - \text{pen}(m')].$$

The corresponding test between m and m' decides m if $T_{m,m'} < 0$, m' if $T_{m,m'} > 0$ and at random if $T_{m,m'} = 0$. Note that the tests corresponding to $T_{m,m'}$ and $T_{m',m}$ are the same. We then set, for all $m \in \mathcal{M}$,

$$\mathcal{R}_m = \{m' \in \mathcal{M}, m' \neq m \mid \text{the test based on } T_{m,m'} \text{ rejects } m\}$$

and, given some $\varepsilon > 0$, we define \hat{m} to be any point in \mathcal{M} such that

$$(5) \quad \mathcal{D}(\hat{m}) \leq \inf_{m \in \mathcal{M}} \mathcal{D}(m) + \varepsilon/3 \quad \text{with} \quad \mathcal{D}(m) = \sup_{m' \in \mathcal{R}_m} \{H^2(\hat{s}_m, \hat{s}_{m'})\}.$$

This model selection procedure results in an estimator $\tilde{s} = \hat{s}_{\hat{m}}$ that we shall call *penalized histogram estimator* (in the sequel PHE, for short) based on the family of models $\{S_m, m \in \mathcal{M}\}$ and the penalty function $\text{pen}(\cdot)$. As to the penalty, it is the sum of two components: $\text{pen}(m) = c_1 |m| + c_2 \Delta_m$ with c_1 and c_2 depending on the framework and Δ_m being a nonnegative weight associated to the model S_m . We require that those weights satisfy

$$(6) \quad \sum_{m \in \mathcal{M}} \exp[-\Delta_m] = \Sigma < +\infty.$$

If $\Sigma = 1$, the choice of the Δ_m can be viewed as the choice of a prior distribution on the models. For related conditions and their interpretation, see Barron and Cover (1991), Barron, Birgé and Massart (1999) or Birgé and Massart (2001). The constant 16 in (4) plays no particular role and has only been chosen in order to improve the legibility of our main results. Our selection procedure can be viewed as a mixture between a method due to Birgé (1983 and 2006) based on testing and

an improved version of the original Lepski's method, as described in Lepski (1991) and subsequent work of the same author. This improved version was presented by Lepski in a series of lectures he gave at Garchy in 1998.

2.4. Risk bounds for the procedure. As we shall see later, with a suitable choice of ε , the performances of this procedure for Problems 0, 1 and 2 are described by risk bounds of the following form:

$$(7) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq C'_0 \inf_{m \in \mathcal{M}} \left\{ \mathbb{E} [(H^2(s, S_m))] + C_P |m| [1 + |m|^{-1} (\Delta_m + \Sigma^2)] \right\},$$

where C'_0 is a numerical constants and C_P as in (3). Comparing (7) with (3), we see that the estimator \tilde{s} achieves a risk bound comparable, up to a constant factor, with the best risk bound obtained by the estimators \hat{s}_m provided that Σ is not large and Δ_m not much larger than $|m|$. Note that these two restrictions are, to some extent, contradictory since the smaller Δ_m , the larger Σ , although it is clearly unnecessary to choose Δ_m smaller than $|m|$. Therefore, if $\sum_{m \in \mathcal{M}} e^{-|m|}$ is not large, one can merely take $\Delta_m = |m|$. Otherwise, the choice of the Δ_m will be more delicate but we should keep in mind that, if Σ is not large, the performance of \tilde{s} will be as good (up to a constant factor) as the performance of any \hat{s}_m for which $\Delta_m \leq |m|$.

3. A REVIEW OF THE MODELS WE SHALL USE

3.1. Some classes of functions of special interest. The motivations for the choice of some family of models $\{S_m, m \in \mathcal{M}\}$ are twofold. First, there is the restriction that \mathcal{M} should satisfy Assumption **H** and there are two main examples of such families. In the "nested" case, the family is totally ordered for the inclusion and thus, we either have $m \vee m' = m$ or $m \vee m' = m'$ for all m and m' in \mathcal{M} . Then, $\overline{\mathcal{M}} = \mathcal{M}$ and $\delta = 1$. Another situation where Assumption **H** is satisfied with $\delta = 1$ occurs when \mathcal{X} is either \mathbb{R} or some subinterval of \mathbb{R} and each $m \in \mathcal{M}$ is a finite partition of \mathcal{X} into intervals.

The second motivation is connected to the approximation properties of the models. If, for instance, we believe that the true s is smooth or monotone, one should introduce families of models that approximate reasonably well such functions. In the sequel, we shall put a special emphasis on the following classes of functions:

- *Monotone functions.* For \mathcal{X} an interval of \mathbb{R} with interior $\overset{\circ}{\mathcal{X}}$ and R a positive number, we denote by $\mathcal{S}^1(R)$ the set of monotone functions t on \mathcal{X} such that $\sup_{x, y \in \overset{\circ}{\mathcal{X}}} |t(x) - t(y)| \leq R$.
- *Continuous functions.* Let w be a modulus of continuity on $[0, 1)$, i.e. a continuous nondecreasing function with $w(0) = 0$ — see additional details in DeVore and Lorentz (1993) —. We denote by $\mathcal{S}^2(w)$ the set of functions t on $[0, 1)$ such that $|t(x+y) - t(x)| \leq w(y)$ for all $x \in [0, 1)$ and $0 \leq y \leq 1-x$. For $0 < \alpha \leq 1$ and $R > 0$, the Hölder class \mathcal{H}_α^R is the class $\mathcal{S}^2(w)$ with $w(y) = Ry^\alpha$. More generally we say that a function u defined on $\mathcal{V} \subset [0, 1)^k$ for some $k \geq 1$ belongs to the set $\mathcal{H}_\alpha^R(\mathcal{V})$, $\alpha \in]0, 1)$, $R > 0$, if

$$|u(x) - u(y)| \leq R \sum_{j=1}^k |x_j - y_j|^\alpha \quad \text{for all } x, y \in \mathcal{V}.$$

- *Piecewise constant functions.* If the function t defined on $[0, 1)$ is constant over some intervals and then jumps from time to time, it is a piecewise constant function of the form

$$(8) \quad t = \sum_{k=1}^D t_k \mathbb{1}_{[x_{k-1}, x_k)} \quad \text{with } 0 = x_0 < x_1 < \dots < x_D = 1.$$

We shall denote by $\mathcal{S}^3(D, R)$ the class of such piecewise functions such that $\sup_{1 \leq k \leq D} t_k \leq R$. Note that this would correspond to a parametric model with D parameters if the locations of the jumps were known. We shall restrict our attention to $D \geq 2$ since $\mathcal{S}^3(1, R)$ only contains constant functions and is then a subset of $\mathcal{S}^2(w)$ with $w \equiv 0$.

- *Besov balls and functions of bounded variation.* Here we consider functions t defined on $[0, 1)$. Given positive numbers α, p and R , we denote by $\mathcal{B}_{p,\infty}^\alpha(R)$, the closed Besov ball of radius R centered at zero of the Besov space $B_{p,\infty}^\alpha([0, 1))$, i.e. the set of functions t in this space with Besov semi-norm $|t|_{B_{p,\infty}^\alpha} \leq R$. Analogously, we set $\mathcal{B}_{BV}(R)$ for the set of functions t of bounded variation with $\text{Var}^*(t) \leq R$. We refer to Chapter 2 of the book by DeVore and Lorentz (1993) for details on Besov spaces and the definition of Besov semi-norms, functions of bounded variation and the variation semi-norm Var^* . Note that $\mathcal{S}_1(R) \subset \mathcal{B}_{BV}(R)$. We shall also consider the multidimensional Besov spaces $B_{p,\infty}^\alpha([0, 1)^k)$ for $k \geq 2$.

3.2. Some typical models. Let us now describe a few useful families of models and corresponding choices for the weights Δ_m that satisfy (6).

3.2.1. Example 1: models for functions on $[0, 1)$. The following models are suitable for approximating functions belonging to the classes that we just mentioned. Since they are based on partitions of $[0, 1)$ into intervals, they satisfy Assumption **H** with $\delta = 1$. Let $\mathcal{J}_l = \{j2^{-l}, j \in \mathbb{N}\}$ and $\mathcal{J}_\infty = \cup_{l \in \mathbb{N}} \mathcal{J}_l$ be the set of all dyadic points in $[0, 1)$. To build \mathcal{M} , we consider partitions $m = \{I_1, \dots, I_D\}$ of $[0, 1)$ generated by increasing sequences $\{0 = x_0 < x_1 < \dots < x_D = 1\}$ with $I_i = [x_{i-1}, x_i)$. We then define \mathcal{M} to be the set of all such partitions with $x_i \in \mathcal{J}_\infty$ for $1 \leq i \leq D-1$. Therefore, whatever $m \in \mathcal{M}$, the elements of \mathcal{S}_m are piecewise constant functions with D pieces and jumps located on the grid \mathcal{J}_∞ . The novelty of this particular family of partitions lies in the fact that there is no lower bound on the length of the intervals on which the partitions are built. It will be useful to single out the set $\mathcal{M}_R = \{m_k, k \in \mathbb{N}\}$ of regular dyadic partitions where m_k is the partition of $[0, 1)$ into 2^k intervals of length 2^{-k} . In particular, $m_0 = [0, 1)$.

One possible way of defining the corresponding weights Δ_m is as follows. For $l \in \mathbb{N}^*$ and $2 \leq D \leq 2^l$ we define $\mathcal{M}_{l,D}$ as the set of all partitions m with $|m| = D$ and l is the smallest integer such that $\{x_1, \dots, x_{D-1}\} \subset \mathcal{J}_l$. Then, $\mathcal{M} = \left[\bigcup_{l \geq 1} \left(\bigcup_{D=2}^{2^l} \mathcal{M}_{l,D} \right) \right] \cup \{m_0\}$. We choose $\Delta_{m_0} = 1$ and

$$(9) \quad \Delta_m = D(l \log 2 + 2 - \log D) + 2 \log l \quad \text{if } m \in \mathcal{M}_{l,D}.$$

Since $|\mathcal{M}_{l,D}| \leq \binom{2^l-1}{D-1} \leq \binom{2^l}{D} \leq (2^l e/D)^D$, we derive from (9) that

$$\begin{aligned} \sum_{m \in \mathcal{M} \setminus \{m_0\}} \exp[-\Delta_m] &< \sum_{l \geq 1} \sum_{D=2}^{2^l} |\mathcal{M}_{l,D}| l^{-2} \exp[-D(l \log 2 + 2 - \log D)] \\ &\leq \sum_{l \geq 1} \sum_{D \geq 2} l^{-2} e^{-D} = \frac{\pi^2 - 6}{6e(e-1)} < 0.14 \end{aligned}$$

and it follows that (6) is satisfied.

3.2.2. Special partitions derived from adaptive approximation algorithms. It is easily seen that the family \mathcal{M} of partitions we introduced for Example 1 is too rich for choosing $\Delta_m = c|m|$ for all m and c a fixed constant since then (6) would not be satisfied. For partitions in $\mathcal{M}_{l,D}$ with $l > D$, Δ_m behaves as $l|m|$ and l can be arbitrarily large. Fortunately, there exists a subset \mathcal{M}_T^1 of \mathcal{M} , which is of special interest because of its approximation properties with respect to functions in Besov spaces, and such as it is possible to choose $\Delta_m = 2|m|$ for $m \in \mathcal{M}_T^1$. This will definitely improve the performances of the PHE for estimating functions in Besov spaces. Let us now describe \mathcal{M}_T^1 .

Among all partitions on $[0, 1)$ with dyadic endpoints, some of them, which are in one-to-one correspondance with the family of complete binary trees, can be derived by the following algorithm described in Section 3.3 of DeVore (1998). One starts with the root of the tree which corresponds to the interval $[0, 1)$ and decides to divide it into two intervals of length $1/2$ or not. We assume here that all intervals contain their left endpoint but not the right one. If one does not divide, the algorithm stops and the tree is reduced to its root. If one divides, one gets two intervals corresponding to adding two sons to the root. Then one repeats the procedure with each interval and so on. . . . At each step, the terminal nodes of the tree correspond to the intervals in the partition and one decides to divide any such interval into two equal parts or not. Dividing means adding two sons to the corresponding terminal node. The whole procedure stops at some stage producing a complete binary tree with D terminal nodes and the corresponding partition of $[0, 1)$ into D intervals. This is the type of tree which comes out of an algorithm like CART, as described by Breiman et al. (1984). Such constructions and the corresponding selection procedure resulting from the CART algorithm have been studied by Gey and Nédélec (2005). We denote by \mathcal{M}_T^1 the subset of \mathcal{M} of all partitions that can be obtained in this way. Note here that the set \mathcal{M}_R of regular partitions is a subset of \mathcal{M}_T^1 .

It is known that the number of complete binary trees with $j+1$ terminal nodes is given by the so-called Catalan numbers $(1+j)^{-1} \binom{2j}{j}$ as explained for instance in Stanley (1999, page 172). As a consequence, we can redefine $\Delta_m = 2|m|$ for $m \in \mathcal{M}_T^1$.

and, using the fact (which derives from Stirling's expansion) that $\binom{2j}{j} \leq 4^j$, get

$$\begin{aligned} \sum_{m \in \mathcal{M}_T^1} \exp[-\Delta_m] &< \sum_{j \geq 0} \sum_{\{m \in \mathcal{M}_T^1 \mid |m|=1+j\}} \exp[-2(j+1)] \\ &= \sum_{j \geq 0} \frac{\binom{2j}{j} \exp[-2(j+1)]}{j+1} \leq e^{-2} \sum_{j \geq 0} \frac{(2/e)^{2j}}{j+1} = \Sigma'_1. \end{aligned}$$

Finally (6) is satisfied with $\Sigma < \Sigma'_1 + 0.14$.

3.2.3. Example 2: estimating functions with radial symmetry. There are situations where one may assume that the value of $s(x)$ only depends on the Euclidean distance $\|x\|$ between this point and some origin in which case one can write $s(x) = \Phi(\|x\|)$. In such a case, it is natural to estimate s on a ball, which we may assume, without loss of generality, to be the open unit ball \mathcal{B}_k of \mathbb{R}^k . To any partition m of $[0, 1)$ we can associate a partition of \mathcal{B}_k with elements $J = \{x \mid \|x\| \in I\}$ where I denotes an element of m . For simplicity, we shall identify the two partitions (the first one of $[0, 1)$ and the new one of \mathcal{B}_k) and denote both of them by m . In the sequel, we shall focus our attention on the family of partitions of Example 1 with the weights defined in Section 3.2.2.

3.2.4. Example 3: estimating functions on $[0, 1)^k$, $k \geq 2$. To deal with the case $\mathcal{X} = [0, 1)^k$, let us first introduce some notations. For $j \in \mathbb{N}$ we consider the set

$$\mathcal{N}_j = \left\{ \mathbf{l} = (l_1, \dots, l_k) \in \mathbb{N}^k \mid 1 \leq l_i \leq 2^j \text{ for } 1 \leq i \leq k \right\}$$

and for $j \in \mathbb{N}$ and $\mathbf{l} \in \mathcal{N}_j$ the cube $K_{j, \mathbf{l}}$ given by

$$K_{j, \mathbf{l}} = \left\{ \mathbf{x} = (x_1, \dots, x_k) \in [0, 1)^k \mid (l_i - 1)2^{-j} \leq x_i < l_i 2^{-j} \text{ for } 1 \leq i \leq k \right\}.$$

We set $\mathcal{K}_j = \{K_{j, \mathbf{l}}, \mathbf{l} \in \mathcal{N}_j\}$ and $\mathcal{K} = \bigcup_{j \geq 0} \mathcal{K}_j$.

Let \mathcal{P} be the collection of all finite subsets p of $\mathcal{K} \setminus \mathcal{K}_0$ consisting of disjoint cubes. To each $p \in \mathcal{P}$, we associate the positive quantity $J(p) = \inf\{j \mid p \cap \mathcal{K}_j \neq \emptyset\}$ ($J(\emptyset) = +\infty$) and the partition m_p generated by p , i.e. $m_p = \{I \in p\} \cup \{[0, 1)^k \setminus \bigcup_{I \in p} I\}$ provided that this last set is not empty and $m_p = \{I \in p\}$ otherwise. We finally set $\mathcal{M} = \{m_p \vee \mathcal{K}_j \text{ with } p \in \mathcal{P} \text{ and } j < J(p)\}$. Note here that the mapping $(j, p) \mapsto m_p \vee \mathcal{K}_j$ is not one to one. For instance $m_\emptyset \vee \mathcal{K}_j = \mathcal{K}_j = \mathcal{K}_j \vee \mathcal{K}_{j-1}$. We shall prove in Section 7.1 the following result:

Lemma 1. *The family \mathcal{M} satisfies Assumption **H** with $\delta = 2$.*

In order to define the weights Δ_m , we shall distinguish a special subset \mathcal{M}_T^k of \mathcal{M} which is the k -dimensional analogue of the one we considered in Section 3.2.2. Here one starts the algorithm with $\mathcal{X} = [0, 1)^k$ (which corresponds to the root of the tree) and at each step get a partition of \mathcal{X} into a finite family of disjoint cubes of the form $K_{j, \mathbf{l}}$. One then decides to divide any such cube into the 2^k elements of \mathcal{K}_{j+1} which are contained in it or not. Again, this corresponds to growing a complete 2^k -ary tree, partitioning a cube meaning adding 2^k sons to a terminal node and the set \mathcal{M}_T^k of all partitions that can be constructed in this way corresponds to the set of complete

2^k -ary trees. As for $k = 1$, \mathcal{M}_T^k contains the set $\mathcal{M}_R = \{m_\emptyset \vee \mathcal{K}_j, j \geq 0\}$ of all regular partitions of \mathcal{X} into 2^{kj} cubes of equal volume. Working with \mathcal{M} instead of the much simpler family \mathcal{M}_R allows to handle less regular functions like those which have a few spikes or are less smooth on some subset of \mathcal{X} .

If $m \in \mathcal{M}_T^k$ we take $\Delta_m = |m|$ and otherwise we set

$$\Delta'_{j,p} = j + k \sum_{i \geq 1} (j+i) |p \cap \mathcal{K}_{j+i}| \quad \text{for } p \in \mathcal{P} \text{ and } j < J(p)$$

and

$$(10) \quad \Delta_m = \inf_{\{(j,p) \mid m = m_p \vee \mathcal{K}_j\}} \{\Delta'_{j,p}\} \quad \text{for } m \in \mathcal{M} \setminus \mathcal{M}_T^k.$$

Note that the ratio $\Delta_m/|m|$ is unbounded for $m \notin \mathcal{M}_T^k$ as shown by the example of $m = m_p \vee \mathcal{K}_0$ with p reduced to a single element of \mathcal{K}_j , $j > 0$. Then $|m| = 2$ while $\Delta_m = kj$ may be arbitrarily large. For the partitions m belonging to \mathcal{M}_T^k we use the fact — see Stanley (1999) — that any complete l -ary tree has a number of terminal nodes of the form $1 + j(l-1)$ for some $j \in \mathbb{N}$ and that the number of such trees with $1 + j(l-1)$ terminal nodes is $[1 + j(l-1)]^{-1} \binom{lj}{j}$. For $l = 2^k$ we derive that the number of partitions in \mathcal{M}_T^k with $1 + j(2^k - 1)$ elements is $[1 + j(2^k - 1)]^{-1} \binom{2^k j}{j}$. Moreover, since $k \geq 2$, we check that

$$\Delta_m > j(k \log 2 + 1) + \log(j+1) \quad \text{if } |m| = 1 + j(2^k - 1).$$

Since $\binom{lj}{j} \leq (le)^j$, it follows that

$$\begin{aligned} \sum_{m \in \mathcal{M}_T^k} \exp[-\Delta_m] &< \sum_{j \geq 0} \sum_{\{m \in \mathcal{M}_T^k \mid |m| = 1 + j(2^k - 1)\}} \frac{\exp[-j(k \log 2 + 1)]}{j+1} \\ &= \sum_{j \geq 0} \frac{\binom{2^k j}{j} (2^k e)^{-j}}{(j+1)[1 + j(2^k - 1)]} \\ &\leq \sum_{j \geq 0} \frac{1}{(j+1)[1 + j(2^k - 1)]} = \Sigma'_k. \end{aligned}$$

Let us now turn to the partitions of the form $m_p \vee \mathcal{K}_j$. For such a partition $p \cap \mathcal{K}_{j'} = \emptyset$ for $j' \leq j$ and, for $i \geq 1$, $|p \cap \mathcal{K}_{j+i}| = l_i$ with $0 \leq l_i \leq 2^{k(j+i)}$. Moreover, the number of those $p \in \mathcal{P}$ such that $|p \cap \mathcal{K}_{j+i}| = l_i$ for a given sequence $\mathbf{l} = (l_i)_{i \geq 1}$ with a finite number of nonzero coefficients is bounded by $\prod_{i \geq 1} \binom{2^{k(j+i)}}{l_i}$. It follows from (10)

that

$$\begin{aligned}
\sum_{m \in \mathcal{M}'} \exp[-\Delta_m] &\leq \sum_{j \geq 0} \sum_{\{p \in \mathcal{P} \mid J(p) > j\}} e^{-j} \prod_{i \geq 1} e^{-k(j+i)|p \cap \mathcal{K}_{j+i}|} \\
&\leq \sum_{j \geq 0} e^{-j} \sum_{\mathbf{l}} \sum_{\{p \mid |p \cap \mathcal{K}_{j+i}| = l_i \text{ for } i \geq 1\}} \prod_{i \geq 1} e^{-k(j+i)l_i} \\
&\leq \sum_{j \geq 0} e^{-j} \sum_{\mathbf{l}} \prod_{i \geq 1} \binom{2^{k(j+i)}}{l_i} e^{-k(j+i)l_i} \\
&\leq \sum_{j \geq 0} e^{-j} \prod_{i \geq 1} \sum_{l_i=0}^{2^{k(j+i)}} \binom{2^{k(j+i)}}{l_i} e^{-k(j+i)l_i} \\
&= \sum_{j \geq 0} e^{-j} \prod_{i \geq 1} \left(1 + e^{-k(j+i)}\right)^{2^{k(j+i)}} \\
&= \sum_{j \geq 0} \exp \left[-j + \sum_{i \geq 1} 2^{k(j+i)} \log \left(1 + e^{-k(j+i)}\right) \right] \\
&\leq \sum_{j \geq 0} \exp \left[-j + \sum_{i \geq 1} (e/2)^{-k(j+i)} \right] = \Sigma_k'' < +\infty.
\end{aligned}$$

Finally we can conclude that (6) holds with $\Sigma < \Sigma'_k + \Sigma_k''$.

3.2.5. Models for n -dimensional vectors. To handle the problem we started with in the introduction, we may assume that our finite index set \mathcal{I} is actually $\mathcal{X} = \{1, \dots, n\}$, the estimation of the function s from \mathcal{X} to \mathbb{R}_+ amounting to the estimation of the vector $(s_1, \dots, s_n)^t \in \mathbb{R}_+^n$ with coordinates $s_i = s(i)$.

Example 4. If one assumes that either s_i varies smoothly with i or is monotone or piecewise constant with a small number of jumps, it is natural to choose for m a partition of \mathcal{X} into intervals and for \mathcal{M} the set of all such partitions. Note that this family satisfies Assumption **H** with $\delta = 1$. Setting here $\Delta_m = |m| + \log \binom{n-1}{|m|-1}$, we get (6) with $\Sigma < (e-1)^{-1}$ since there are $\binom{n-1}{D-1}$ partitions in \mathcal{M} with D elements for $1 \leq D \leq n$.

Example 5. An alternative case is the case when s is constant, equal to \bar{s} on \mathcal{X} except for a few number of locations i where $s(i) \neq \bar{s}$. Since the number k of such locations is unknown, it is natural, for each $k \in \{0, \dots, n-1\}$ to define \mathcal{M}_k as the set of partitions of \mathcal{X} with k singletons and the set of the $n-k$ remaining points. We finally set $\mathcal{M} = \cup_{0 \leq k \leq n-1} \mathcal{M}_k$. Then Assumption **H** holds with $\delta = 1$. For $m \in \mathcal{M}_k$, $|m| = k+1$ and we set $\Delta_m = \log \binom{n}{k} + k = \log \binom{n}{|m|-1} + |m| - 1$, so that (6) holds with $\Sigma < e/(e-1)$.

4. THE CASE OF A DETERMINISTIC MEASURE M

Let us now see how our general framework applies to Problems 1 and 2. Besides these, our setup also covers the problem of density estimation. Although there is a huge amount of literature on density estimation, our method brings some improvements to known results on partition selection for histograms. Moreover, since this problem has attracted so much attention, it can serve as pedagogical example and reference for the sequel. This is why, before considering more original and less studied frameworks, we shall start our review by this quite familiar estimation problem.

4.1. Density estimation. We consider the classical problem of estimating an unknown density s from a sample of size n , which means that we have at hand an i.i.d. sample X_1, \dots, X_n from a distribution with unknown density s with respect to some given measure $M = \lambda$ on \mathcal{X} . We define N to be the empirical distribution: $N(A) = n^{-1} \sum_{i=1}^n \mathbb{1}_{X_i \in A}$. Then, as required, $\mathbb{E}[N(A)] = \int_A s d\lambda$ for all measurable subsets A of \mathcal{X} . In this case the distance H is merely a version of the Hellinger distance between densities.

Within this framework, we can prove the following general result.

Theorem 1. *Assume that the family \mathcal{M} satisfies Assumption **H** and the weights $\{\Delta_m, m \in \mathcal{M}\}$ are chosen so that (6) holds. Then the penalized histogram estimator $\tilde{s} = \hat{s}_{\hat{m}}$ defined in Section 2.3 with $\text{pen}(m) \geq n^{-1}(8\delta|m| + 202\Delta_m)$ satisfies*

$$\mathbb{E}[H^2(\tilde{s}, s)] \leq \left[390 \left(\inf_{m \in \mathcal{M}} (H^2(s, S_m) + \text{pen}(m)) + \frac{101\Sigma^2}{n} \right) + \varepsilon \right] \wedge 2.$$

The only previous works on partition selection for histograms using squared Hellinger loss we know about are to be found in Castellan (1999 and 2000) and Birgé (2006). Castellan's approach is based on penalized maximum likelihood. This requires to make specific restrictions on the underlying density s , in particular that s should be bounded away from 0. For the problem of estimating a density on \mathbb{R} , her conditions on the family of partitions are also more restrictive than ours since we can handle any countable families of finite partitions into intervals. Nevertheless, in the multivariate case, our assumptions on the partitions are more stringent. Birgé's approach based on aggregation of histograms built on one half of the sample leads to more abstract but more general results.

Let us now apply the above theorem to various families of models, systematically setting $\text{pen}(m) = n^{-1}(8\delta|m| + 202\Delta_m)$ and $\varepsilon = n^{-1}$. We assume in this section that λ is the Lebesgue measure on \mathcal{X} .

4.1.1. Example 1, continued. When $\mathcal{X} = [0, 1)$, we use the family of models and weights of Section 3.2.1. Our next proposition shows that the PHE based on this simple family of models and weights has nice properties for estimating various types of functions. The proof will be given in Section 7.3.

Proposition 1. *Let \tilde{s} be the PHE based on the family of models and weights Δ_m defined in Section 3.2.1, $\varepsilon = n^{-1}$ and the penalty function $\text{pen}(m) = n^{-1}(8\delta|m| + 202\Delta_m)$.*

i) If $s \in \mathcal{S}_1(R)$, then

$$(11) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq C \left\{ [Rn^{-1} \log(1 + nR^2)]^{2/3} \vee n^{-1} \right\}.$$

ii) If $\sqrt{s} \in \mathcal{S}^2(w)$ where w is a modulus of continuity on $[0, 1]$, we define x_w to be the unique solution of the equation $n x w^2(x) = 1$ if $w(1) \geq n^{-1/2}$ and $x_w = 1$ otherwise. Then

$$(12) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq C(n x_w)^{-1}.$$

If, in particular, \sqrt{s} belongs to the Hölder class \mathcal{H}_α^R with $R \geq n^{-1/2}$, then $\mathbb{E}_s [H^2(\tilde{s}, s)] \leq C R^{2/(2\alpha+1)} n^{-2\alpha/(2\alpha+1)}$.

iii) If $s \in \mathcal{S}^3(D, R)$ with $2 \leq D \leq n$ and $R \geq 2$, we get

$$(13) \quad \mathbb{E}_s [H^2(\tilde{s}, s)] \leq C D n^{-1} \log(nR/D).$$

It is interesting to see to what extent the previous bounds (together with the trivial one, $\mathbb{E} [H^2(\tilde{s}, s)] \leq 2$, which always holds but which we did not include in (11), (12) and (13) for simplicity) are optimal (up to the universal constants C). Many lower bounds on the minimax risk over various density classes are known for classical loss functions. For squared Hellinger loss, some are given in Birgé (1983 and 1986) and Birgé and Massart (1998). Many more are known for the squared \mathbb{L}_2 -loss, which can easily be extended to squared Hellinger loss because their proofs are based on perturbations arguments involving sets of densities for which both distances are equivalent. It follows from these classical results that the bound we find for continuous densities are actually optimal (see Birgé, 1983, p.211) while (11) is suboptimal because of the presence of the log factor. We shall see below that the more sophisticated penalization strategy introduced in Section 3.2.2 does solve the problem. The case of piecewise constant functions is more complicated. If D and the locations of the jumps were known, one could use a single model corresponding to the relevant partition with D intervals and get a risk bound CD/n corresponding to a parametric problem with D parameters. Apart from the constant C , this bound cannot be improved which shows that the study of uniform risk bounds over $\mathcal{S}^3(D, R)$ is only of interest when $D \leq n$ since otherwise a lower bound for the risk is of the order of the trivial upper bound 2. When D is smaller than n the extra $\log(nR/D)$ factor in (13) is due to the fact that we have to estimate the locations of the jumps. The problem has been considered in Birgé and Massart (1998, Section 4.2 and Proposition 2) where it is shown that a lower bound for the risk (when $n \geq 5D$ and $D \geq 9$) is $c D n^{-1} \log(nD^{-1})$. Therefore our bound is optimal for moderate values of R . We do not know whether the $\log R$ factor in the upper bound is necessary or not.

4.1.2. Improved risk bounds with a better weighting strategy. If we use the weights Δ_m defined in Section 3.2.2 to build \tilde{s} , we can only improve (up to constants) the risk bounds given in Proposition 1 since the value of Σ does not change much while the new weights are not larger than the previous ones. Besides, the values of the weights have been substantially decreased for the partitions belonging to \mathcal{M}_T^1 . It turns out that piecewise constants functions on the elements of \mathcal{M}_T^1 possess quite powerful approximation properties with respect to functions in Besov spaces $B_{p,\infty}^\alpha([0, 1])$ with $\alpha < 1$ and monotone functions. These properties are given in the following theorem which also includes the multidimensional case.

Theorem 2. Let $\mathcal{X} = [0, 1]^k$, \mathcal{M}_T^k be the set of partitions m of \mathcal{X} defined in Section 3.2.4 and, for $m \in \mathcal{M}_T^k$, let S'_m be the cone $\{t = \sum_{I \in m} t_I \mathbb{1}_I, t_I \geq 0\}$. For any $p > 0$, α with $1 > \alpha > k(1/p - 1/2)_+$ and any function t belonging to the Besov space $B_{p,\infty}^\alpha([0, 1]^k)$ with Besov semi-norm $|t|_{B_{p,\infty}^\alpha}$, one can find some $t' \in \bigcup_{m \in \mathcal{M}_T^k} S'_m$ such that

$$(14) \quad \|t - t'\|_2 \leq C(\alpha, k, p) |t|_{B_{p,\infty}^\alpha} |m|^{-\alpha/k},$$

where $\|\cdot\|_2$ denotes the $\mathbb{L}_2(dx)$ -norm on $[0, 1]^k$.

If t is a function of bounded variation on $[0, 1]$, there exists $t' \in \bigcup_{m \in \mathcal{M}_T^1} S'_m$ such that $\|t - t'\|_2 \leq C' \text{Var}^*(t) |m|^{-1}$.

The bound (14) is given in DeVore and Yu (1990). The proof for the bounded variation case has been kindly communicated to the second author by Ron DeVore. With the help of this theorem, we can now derive from Theorem 1 the following improved bounds the proof of which is straightforward.

Proposition 2. Let \tilde{s} be the PHE based on the weights Δ_m defined in Section 3.2.2. If \sqrt{s} is a function of bounded variation with $\text{Var}^*(\sqrt{s}) \leq R$ and in particular if it belongs to $\mathcal{S}^1(R)$, then

$$(15) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq \min \left\{ C(R/n)^{2/3}, 2 \right\} \quad \text{for } R \geq n^{-1/2}.$$

If $\sqrt{s} \in B_{p,\infty}^\alpha([0, 1])$ with $1 > \alpha > (1/p - 1/2)_+$ and $|\sqrt{s}|_{B_{p,\infty}^\alpha} \leq R$ with $R \geq n^{-1/2}$, then

$$\mathbb{E} [H^2(\tilde{s}, s)] \leq \min \left\{ C(\alpha, p) R^{2/(1+2\alpha)} n^{-2\alpha/(1+2\alpha)}, 2 \right\}.$$

It follows from classical lower bounds arguments that these bounds are minimax up to constants.

4.1.3. The multidimensional case. When the density s defined on $\mathcal{X} = \mathcal{B}_k$ can be written $s(x) = \Phi(\|x\|)$ for some function Φ on $[0, 1]$, we use the family of models introduced in Example 2. We then obtain the risk bounds given in Propositions 1 and 2 if we replace the assumptions on s by the same on Φ . We omit the details.

If $\mathcal{X} = [0, 1]^k$, $k \geq 2$ and we use the family of models and weights described in Section 3.2.4, we get the following result.

Proposition 3. Let $R \geq k^{-1} n^{-1/2}$. If \sqrt{s} belong to $\mathcal{H}_\alpha^R([0, 1]^k)$, then

$$(16) \quad \mathbb{E} [H^2(s, \tilde{s})] \leq \min \left\{ C(Rk)^{2k/(k+2\alpha)} n^{-2\alpha/(2\alpha+k)}, 2 \right\}.$$

More generally, if \sqrt{s} belongs to $B_{p,\infty}^\alpha([0, 1]^k)$ with $1 > \alpha > k(1/p - 1/2)_+$ and $|\sqrt{s}|_{B_{p,\infty}^\alpha} \leq R$, then

$$\mathbb{E} [H^2(\tilde{s}, s)] \leq \min \left\{ C(\alpha, k, p) R^{2k/(k+2\alpha)} n^{-2\alpha/(k+2\alpha)}, 2 \right\}.$$

Proof: Let $m = \mathcal{K}_j$ be an element of \mathcal{M}_R . Then $\Delta_m = |m| = 2^{kj}$ and the maximal variation of a function of $\mathcal{H}_\alpha^R([0, 1]^k)$ on an element of m is bounded by $Rk2^{-j\alpha}$ so that $H^2(s, S_m) \leq (Rk)^2 2^{-2j\alpha}$. It then follows from Theorem 1 that $\mathbb{E} [H^2(\tilde{s}, s)] \leq$

$C' [(Rk)^2 2^{-2j\alpha} + n^{-1} 2^{kj}]$. The lower bound on R allows us to choose $j \in \mathbb{N}$ such that $2^j \leq (n(Rk)^2)^{1/(k+2\alpha)} < 2^{j+1}$ which leads to

$$\mathbb{E} [H^2(\tilde{s}, s)] \leq C' \left[(Rk)^2 2^{2\alpha} (n(Rk)^2)^{-2\alpha/(k+2\alpha)} + n^{-1} (n(Rk)^2)^{k/(k+2\alpha)} \right].$$

The first bound follows since $2^{2\alpha} \leq 4$. The second bound can be proved in the same way from (14). \square

4.2. Poisson processes. Let us consider the stochastic framework corresponding to Problem 2 where ν is dominated by some given measure $M = \lambda$ on \mathcal{X} with density $s = d\nu/d\lambda$. This implies that (1) holds as required. In this case, the performances of the PHE \tilde{s} are as follows.

Theorem 3. *Assume that the family \mathcal{M} satisfies Assumption **H** and the weights $\{\Delta_m, m \in \mathcal{M}\}$ are chosen so that (6) holds. Then the estimator \tilde{s} defined in Section 2.3 with $\text{pen}(m) \geq 3\delta|m| + 6\Delta_m$ satisfies*

$$(17) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq 390 \left[\inf_{m \in \mathcal{M}} [H^2(s, S_m) + \text{pen}(m)] + 3\Sigma^2 \right] + \varepsilon.$$

This theorem should be compared with the results of Reynaud-Bouret (2003) who uses more general families of projection estimators than just histograms based on partitions. Nevertheless, for the problem we consider here, her choice of the \mathbb{L}^2 -loss induces some restrictions on both the intensity and the collection of partitions at hand. For instance, the intensity has to be bounded and the procedure requires some suitable estimation of its sup-norm. As Castellán (1999), she cannot deal with partitions with arbitrary small length.

Let us now apply this theorem to our families of models, systematically setting $\text{pen}(m) = 3\delta|m| + 6\Delta_m$ and $\varepsilon = 1$. In view of facilitating the interpretation of the results to follow, it is convenient to use an analogy with density estimation. This analogy, based on the following heuristics, allows to extrapolate the bounds from one framework to the other.

We recall that observing the Poisson process N of intensity s is equivalent to observing \overline{N} i.i.d. random variables with density s' , where $\overline{N} = N(\mathcal{X})$ is a Poisson variable with parameter $n = \int_{\mathcal{X}} s d\lambda$ and $s' = n^{-1}s$. With this in mind, and even though n need not be an integer, we can view the estimation of s as an analogue of the estimation of the density s' from n i.i.d. observations. Pursuing into this direction, we may rewrite the risk in the Poisson case as $\mathbb{E} [H^2(\tilde{s}, s)] = n\mathbb{E} [H^2(n^{-1}\tilde{s}, s')]$ and, setting $\tilde{s}_n = n^{-1}\tilde{s}$, view $\mathbb{E} [H^2(\tilde{s}_n, s')] = n^{-1}\mathbb{E} [H^2(\tilde{s}, s)]$ as an analogue of the risk for estimating s' from n i.i.d. observations. When \sqrt{s} belongs to $\mathcal{S}^1(R)$, $\mathcal{S}^2(w)$ or $\mathcal{S}^3(D, R)$, then the square-root of the density $s' = s/n$ belongs to $\mathcal{S}^1(Rn^{-1/2})$, $\mathcal{S}^2(wn^{-1/2})$ or $\mathcal{S}^3(D, Rn^{-1/2})$ respectively (provided that $R^2 \geq n$ in the last case, since otherwise $\mathcal{S}^3(D, Rn^{-1/2})$ would not contain any density). From these two remarks, we may conclude that a risk bound of the form $f(R)$ in the Poisson case should be interpreted in the density case as $n^{-1}f(Rn^{-1/2})$.

Example 1, continued. Here we deal with a Poisson process N on a finite interval of \mathbb{R} , which we may assume, without loss of generality, to be $[0, 1]$, of intensity s with respect to the Lebesgue measure ν . To estimate s we use the family of models of Example 1 with the weights Δ_m defined in Section 3.2.2. The resulting

PHE \tilde{s} has the following properties which can be proved exactly like those given in Propositions 1 and 2.

Proposition 4. *Let w be a modulus of continuity on $[0, 1)$. We define x_w to be the unique solution of the equation $xw^2(x) = 1$ if $w(1) \geq 1$ and $x_w = 1$ otherwise. Then*

$$(18) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq Cx_w^{-1} \quad \text{for all } s \text{ such that } \sqrt{s} \in \mathcal{S}^2(w).$$

If, in particular, \sqrt{s} belongs to the Hölder class \mathcal{H}_α^R with $R \geq 1$, then

$$\mathbb{E} [H^2(\tilde{s}, s)] \leq CR^{2/(2\alpha+1)}.$$

Given $D \geq 2$ and $R \geq 2D$, we get

$$(19) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq CD \log(R/D) \quad \text{for all } s \in \mathcal{S}^3(D, R).$$

If \sqrt{s} belongs to $\mathcal{S}^1(R)$ with $R \geq 1$, then $\mathbb{E} [H^2(\tilde{s}, s)] \leq CR^{2/3}$.

If $\sqrt{s} \in B_{p,\infty}^\alpha([0, 1))$ with $1 > \alpha > (1/p - 1/2)_+$ and $|\sqrt{s}|_{B_{p,\infty}^\alpha} \leq R$ with $R \geq 1$, then $\mathbb{E} [H^2(\tilde{s}, s)] \leq CR^{2/(1+2\alpha)}$.

For the sake of simplicity, let us assume that $n = \int_{\mathcal{X}} s d\lambda$ is an integer. The connection established above between the estimation of a density and that of the intensity of a Poisson process shows that Proposition 4 is actually a perfect analogue of Propositions 1 and 2. Namely, when \sqrt{s} belongs to $\mathcal{S}^1(R)$ or $\mathcal{S}^2(w)$ or $s \in \mathcal{S}^3(D, R)$ and $s' = s/n$ then $\sqrt{s'}$ respectively belongs to $\mathcal{S}^1(Rn^{-1/2})$ or $\mathcal{S}^2(wn^{-1/2})$ or $s' \in \mathcal{S}^3(D, Rn^{-1})$ and the risk bounds we get for estimating the intensity s (with respect to the H^2/n -loss) are the same as those obtained from a n sample for estimating the density s' (with the H^2 -loss).

Example 2, continued. If we observe a Poisson process on $\mathcal{X} = \mathcal{B}_k$ with intensity $s(x) = \Phi(\|x\|)$ with respect to the Lebesgue measure for Φ some function on $[0, 1)$ and consider the family of models introduced in Example 1 we obtain the risk bounds given in Proposition 4 if we replace the assumptions on s by the same on Φ .

Example 3, continued. If $\mathcal{X} = [0, 1)^k$ with $k \geq 2$, we use the models and weights defined in Section 3.2.4. Proceeding as for Proposition 3 we get:

Proposition 5. *Let \sqrt{s} belong to $\mathcal{H}_\alpha^R([0, 1)^k)$, then*

$$\mathbb{E}_s [H^2(s, \tilde{s})] \leq C(Rk \vee 1)^{2k/(k+2\alpha)}.$$

If \sqrt{s} belongs to $B_{p,\infty}^\alpha([0, 1)^k)$ with $1 > \alpha > k(1/p - 1/2)_+$ and $|\sqrt{s}|_{B_{p,\infty}^\alpha} \leq R$, then

$$\mathbb{E} [H^2(\tilde{s}, s)] \leq C(\alpha, k, p)(R \vee 1)^{2k/(k+2\alpha)}.$$

As shown by the proof of Proposition 3, we only use the partitions in \mathcal{M}_R to get (16) so that it would be of little use to introduce other partitions if we only wanted to estimate intensities such that \sqrt{s} belong to $\mathcal{H}_\alpha^R([0, 1)^k)$. The interest of considering the larger family \mathcal{M} and to have a special definition of Δ_m when $m \in \mathcal{M}_T$ is that it allows to improve the results when we deal with less regular functions than those for which \sqrt{s} belong to $\mathcal{H}_\alpha^R([0, 1)^k)$, in particular those functions that belong to Besov spaces $B_{p,\infty}^\alpha([0, 1)^k)$ with $1 > \alpha > k/p$. To illustrate this fact, let us study the estimation of those intensities s such that \sqrt{s} has the following specific structure.

Given the nonempty set \mathcal{V} which is a finite union of elements of \mathcal{K} , there is a smallest integer \bar{j} such that \mathcal{V} can be written as the union of N elements of $\mathcal{K}_{\bar{j}}$ with a volume $V = N2^{-k\bar{j}} > 0$. To avoid trivialities, we assume that $\bar{j} > 0$, hence $V < 1$.

Proposition 6. *Let s be an intensity on $[0, 1]^k$ such that $\sqrt{s}\mathbb{1}_{\mathcal{V}}$ belongs to $\mathcal{H}_{\alpha}^R(\mathcal{V})$ with $R \geq 1$ while $\sqrt{s}\mathbb{1}_{\mathcal{V}^c}$ is constant and let \tilde{s} be the PHE based on the weights Δ_m defined in Section 3.2.4. Then*

$$(20) \quad \mathbb{E}[H^2(\tilde{s}, s)] \leq C \inf_{m \in \mathcal{M}} B_m \quad \text{with } B_m = H^2(s, S_m) + |m| + \Delta_m$$

and

$$(21) \quad B_m \leq C \min \left\{ 2^{k\bar{j}} + V^{k/(k+2\alpha)} (kR)^{2k/(k+2\alpha)} ; \right.$$

$$(22) \quad \left. V \left[k\bar{j}2^{k\bar{j}} + (kR)^{2k/(2\alpha+k)} [\log(Rk)]^{2\alpha/(2\alpha+k)} \right] ; \right.$$

$$(23) \quad \left. V \left[2^{k\bar{j}}2^{k\bar{j}} + (kR)^{2k/(2\alpha+k)} \right] \right\}.$$

Proof: Since (20) is merely a consequence of Theorem 3 with the choice $\text{pen}(m) = 3\delta|m| + 6\Delta_m$ and $\varepsilon = 1$, we only have to bound B_m . Let us first consider a regular partition $m = \mathcal{K}_j$. If $j < \bar{j}$, the bias $H^2(s, S_m)$ may be arbitrarily large since the intensity s may be arbitrarily large on \mathcal{V} while it may be small on \mathcal{V}^c . For $j \geq \bar{j}$, the argument used for the proof of Proposition 3 shows that on \mathcal{V} , \sqrt{s} can be approximated uniformly by an element of S_m with a precision at least $Rk2^{-j\alpha}$ so that $H^2(s, S_m) \leq VR^2k^22^{-2j\alpha}$ and $B_m \leq VR^2k^22^{-2j\alpha} + 2^{kj+1}$. If $[VR^2k^2]^{1/(2\alpha+k)} \leq 2^{\bar{j}}$ we set $j = \bar{j}$ and otherwise choose j so that $2^j \leq [VR^2k^2]^{1/(2\alpha+k)} < 2^{j+1}$. This leads to (21).

If we set $m = m_p \vee \mathcal{K}_0$ with p being the set of those $N2^{k(j-\bar{j})} = V2^{kj} \geq 1$ elements of \mathcal{K}_j ($j \geq \bar{j} \geq 1$) that exactly cover \mathcal{V} , we get, since $k \geq 2$

$$B_m \leq VR^2k^22^{-2j\alpha} + (kj + 1)V2^{kj} + 1 \leq Vk \left[R^2k^22^{-2j\alpha} + 2j2^{kj} \right].$$

If $[k^2R^2/\log(kR)]^{1/(2\alpha+k)} < 2^{\bar{j}}$ we set $j = \bar{j}$ and otherwise choose j so that $2^j \leq [k^2R^2/\log(kR)]^{1/(2\alpha+k)} < 2^{j+1}$ which finally leads to (22).

To study the approximation properties of the elements of \mathcal{M}_T^k let us consider a particular cube $K' = K_{\bar{j}, \mathbf{l}} \in \mathcal{V} \cap \mathcal{K}_{\bar{j}}$. Identifying the partitions in \mathcal{M}_T^k with the trees from which they derive, we can design an element $m_{K'}$ of \mathcal{M}_T^k with $2^k - 1$ terminal nodes at each level 1 to \bar{j} and the remaining node K' at level \bar{j} . Then we keep only non-terminal nodes up to level $j \geq \bar{j}$, all nodes at this last level j being terminal, so that their number is $2^{k(j-\bar{j})}$. The total number of terminal nodes of the tree is therefore $\bar{j}(2^k - 1) + 2^{k(j-\bar{j})}$. We can repeat this operation for each of the N cubes in $\mathcal{V} \cap \mathcal{K}_{\bar{j}}$ keeping the value of j fixed. This results in N similar trees. We finally consider the smallest complete tree m that contains the N previous ones. Its number of terminal nodes is then bounded by $N[\bar{j}(2^k - 1) + 2^{k(j-\bar{j})}]$ so that

$$B_m \leq V(Rk)^22^{-2j\alpha} + 2N[\bar{j}(2^k - 1) + 2^{k(j-\bar{j})}] \leq 2V \left[R^2k^22^{-2j\alpha} + \bar{j}2^{k(\bar{j}+1)} + 2^{kj} \right].$$

If $(k^2R^2)^{1/(2\alpha+k)} < 2^{\bar{j}}$ we set $j = \bar{j}$ and otherwise choose j so that $2^j \leq (k^2R^2)^{1/(2\alpha+k)} < 2^{j+1}$, which leads to (23). \square

A comparison of the three bounds (21), (22) and (23) shows that (23) is always

better if we omit the influence of \bar{j} and k but the situation becomes more involved if we take into account the effect of k and \bar{j} . Depending on the values of V, R, \bar{j}, α and k , each type of partition may be the best which justifies to introduce them all.

Remark: An analogue of Proposition 6 holds for density estimation.

4.3. Non-negative random vectors. Let us recall from the introduction that we observe an n -dimensional random vector with independent nonnegative components N_1, \dots, N_n and respective distributions depending on positive parameters s_1, \dots, s_n . One should think of the N_i as Poisson or binomial random variables with unknown expectations s_i . More generally, we assume that there exist some known constants $\kappa > 0$ and $\tau \geq 0$ such that for all $i \in \mathcal{X} = \{1, \dots, n\}$

$$(24) \quad \log \left(\mathbb{E} \left[e^{z(N_i - s_i)} \right] \right) \leq \kappa \frac{z^2 s_i}{2(1 - z\tau)} \quad \text{for all } z \in \left[0, \frac{1}{\tau} \right],$$

with the convention $1/\tau = +\infty$ if $\tau = 0$, and

$$(25) \quad \log \left(\mathbb{E} \left[e^{-z(N_i - s_i)} \right] \right) \leq \kappa \frac{z^2 s_i}{2} \quad \text{for all } z \geq 0.$$

In the case of Poisson or binomial random variables, one can take $\kappa = \tau = 1$ as we shall see below.

Our aim is to estimate the function s from \mathcal{X} to \mathbb{R}_+ given by $s(i) = s_i$. Here we denote by λ the counting measure on \mathcal{X} and set $Y \equiv 1$. Hence $M = \lambda$ and $N(A) = \sum_{i \in A} N_i$. Then \mathcal{L} can be identified with \mathbb{R}_+^n , $\mathbb{E}[N(A)] = \int_A s d\lambda$ as required and $H^2(t, t') = \sum_{i=1}^n \left[\sqrt{t(i)} - \sqrt{t'(i)} \right]^2$ for $t, t' \in \mathcal{L}$.

Theorem 4. Assume that (24) and (25) hold, that the family \mathcal{M} satisfies Assumption **H** and the weights $\{\Delta_m, m \in \mathcal{M}\}$ are chosen so that (6) holds. Let $\text{pen}(m) \geq \kappa [\delta (1 + K^2) |m| + 3K^2 \Delta_m]$ with

$$K = \sqrt{2} \quad \text{if } \tau \leq \kappa; \quad K = \frac{\sqrt{2}}{2} + \sqrt{\frac{\tau}{\kappa} - \frac{1}{2}} \quad \text{if } \tau > \kappa;$$

and let \tilde{s} be the PHE defined in Section 2.3. Then

$$\mathbb{E} [H^2(\tilde{s}, s)] \leq 390 \left[\inf_{m \in \mathcal{M}} [H^2(s, S_m) + \text{pen}(m)] + (3/2) \kappa K^2 \Sigma^2 \right] + \varepsilon.$$

Let us first check that some classical distributions do satisfy Inequalities (24) and (25). If N_i is a binomial random variable with parameters n_i, p_i then for all $z \in \mathbb{R}$,

$$(26) \quad \log \left(\mathbb{E} \left[e^{z(N_i - s_i)} \right] \right) \leq s_i (e^z - z - 1) \quad \text{with } s_i = n_i p_i.$$

If N_i is a Poisson random variable with parameter s_i , then equality holds in (26). Using the bounds $e^z - z - 1 \leq z^2/[2(1 - z)]$ for $z \in [0, 1]$ and $e^z - z - 1 \leq z^2/2$ for $z < 0$ we derive that, in both cases, (24) and (25) hold with $\kappa = \tau = 1$. If N_i has a Gamma distribution $\Gamma(s_i, 1)$, $\mathbb{E}[N_i] = s_i$ and, following the proof of Lemma 1 of Laurent and Massart (2000), we deduce that (24) and (25) hold again with $\kappa = \tau = 1$. More generally, it follows from some version of Bernstein's Inequality — see Lemma 8 of Birgé and Massart (1998) — that (24) holds as soon as

$$\mathbb{E} [(N_i)^p] \leq \kappa \frac{p!}{2} s_i \tau^{p-2}, \quad \text{for all } i \in \mathcal{X} \quad \text{and } p \geq 2.$$

Inequality (25) is always satisfied if $N_i \leq \kappa$. Indeed it follows from

$$e^{-zx} \leq 1 - zx + z^2 x^2 / 2, \quad \forall x, z \geq 0$$

that all non-negative random variables X bounded by κ satisfy

$$\mathbb{E}[e^{-zX}] \leq 1 - z\mathbb{E}[X] + \frac{z^2 \mathbb{E}[X^2]}{2} \leq \exp(-z\mathbb{E}[X] + \kappa z^2 \mathbb{E}[X]/2).$$

The results of Kolaczyk and Nowak (2004), which are based on some sort of discretized penalized maximum likelihood estimator in the spirit of Barron and Cover (1991), have some similarity with ours but they assume that the components of the vector s belong to some known interval $[c, C]$, $c > 0$ and they explicitly use the values of c and C in the construction of their estimator. Such an assumption, which implies, as in the case of density estimation, that squared Hellinger distance and Kullback divergence are equivalent also greatly simplifies the estimation problem.

Example 4, continued. Setting

$$(27) \quad \text{pen}(m) = \kappa [(1 + K^2)|m| + 3K^2 \Delta_m] \quad \text{and} \quad \varepsilon = 1.$$

and using $\log \binom{n-1}{D-1} \leq (D-1)(1 + \log[(n-1)/(D-1)])$ with the convention $0 \log((n-1)/0) = 0$ we get the risk bound

$$(28) \quad \mathbb{E}[H^2(\tilde{s}, s)] \leq C(\kappa, K) \inf_{m \in \mathcal{M}} \left\{ H^2(s, S_m) + |m| + (|m| - 1) \log \left(\frac{n-1}{|m|-1} \right) \right\}.$$

If, for instance, s itself belongs to some S_m with a small value of $|m|$, which corresponds to a piecewise stationary process $(N_i)_{1 \leq i \leq n}$ with a few distribution changes, the risk is bounded by $C(\kappa, K)|m| \log n$.

Another interesting situation corresponds to the case of a monotone sequence $(s_i)_{1 \leq i \leq n}$, i.e. a monotone function s on \mathcal{X} that we may assume, without loss of generality to be nondecreasing.

Proposition 7. *Let the sequence s_i , $1 \leq i \leq n$ be nondecreasing with $\sqrt{s_n} - \sqrt{s_1} = R$, then the PHE \tilde{s} based on the models of Example 4 with pen and ε given by (27) satisfies the following risk bounds with a constant C depending only on κ and K :*

- if $R^2 \leq n^{-1} \log n$, then $\mathbb{E}[H^2(\tilde{s}, s)] \leq C(\kappa, K)(nR^2 + 1)$;
- if $R \geq n/\sqrt{3}$, then $\mathbb{E}[H^2(\tilde{s}, s)] \leq C(\kappa, K)n$;
- otherwise $\mathbb{E}[H^2(\tilde{s}, s)] \leq C(\kappa, K)[R\sqrt{n} \log(n/R)]^{2/3}$.

Remark: If we restrict ourselves to the case $n = 2^k$, we can turn any function s on \mathcal{X} into a function s' on $[0, 1)$ by setting $s' = \sum_{i=1}^n s(i) \mathbb{1}_{[(i-1)2^{-k}, i2^{-k})}$. This transformation will, in particular, preserve the monotonicity properties of the functions. One could then estimate s' using the more sophisticated families of weights that we introduced in Section 3.2.2. The use of this strategy would improve the estimation of monotone functions, removing the logarithmic factors.

Example 5, continued. Choosing pen and ε as in (27) and using the same arguments as for Example 1, we derive an analogue of (28) with n replacing $n-1$ in the logarithmic factor. If we assume that $s_i = \bar{s}$ for $i \notin I$ with $|I| = k$, then $H^2(s, S_m) = 0$ for some $m \in \mathcal{M}_k$ and

$$\mathbb{E}[H^2(\tilde{s}, s)] \leq C(\kappa, K)[k + 1 + k \log(n/k)].$$

5. SPECIAL COUNTING PROCESSES ON THE LINE

Let \mathcal{X} be some interval of \mathbb{R}_+ of the form $[0, \zeta)$ where $0 < \zeta \leq +\infty$ with its Borel σ -algebra \mathcal{A} . We recall that a (univariate) counting process \tilde{N} on \mathcal{X} is a cadlag (right-hand continuous and left-hand limited) process from \mathcal{X} to \mathbb{R}_+ , vanishing at time $t = 0$, with piecewise constant and nondecreasing paths having jumps of size +1 only. The use of counting processes in statistical modeling is developed in great details in the book by Andersen *et al.* (1993) where the interested reader will find many concrete situations for which these processes naturally arise. Typically, \tilde{N}_t counts the number of occurrences of a certain event from time 0 up to time t . The *jumping times* of the process give the dates of occurrence of the event. A counting process can be associated to a random measure N on \mathcal{X} whose cumulative distribution function is the counting process itself, i.e. $N([0, t]) = \tilde{N}_t$ for all $t \in \mathcal{X}$. In the sequel, we shall not distinguish between the counting process \tilde{N} and its associated measure N .

In this paper, we consider a phenomenon which is described by some bounded counting process N^* on \mathcal{X} such that $N^*(\mathcal{X}) \leq k$ a.s. for some known integer k . This means that N^* describes an event that occurs at most k times during the period \mathcal{X} . We also assume that there exist a deterministic measure λ on \mathcal{X} , a deterministic nonnegative function $s \in \mathbb{L}_1(\mathcal{X}, d\lambda)$ and a nonnegative observable process Y^* bounded by 1 on \mathcal{X} such that

$$(29) \quad \mathbb{E}[N^*([0, t])] = \mathbb{E}\left[\int_0^t s Y^* d\lambda\right] \quad \text{for all } t \in \mathcal{X}.$$

We actually observe an aggregated counting process N which is the sum of n i.i.d. processes N^j , $j = 1, \dots, n$ with the same distribution as N^* . The fact that the measure N^j is determined by its cumulative distribution function and (29) imply that there are i.i.d. observable processes Y^j , $j \in \{1, \dots, n\}$ with the distribution of Y^* such that

$$\mathbb{E}[N^j(A)] = \mathbb{E}\left[\int_A s Y^j d\lambda\right] \quad \text{for all } A \in \mathcal{A} \quad \text{and} \quad 1 \leq j \leq n.$$

Therefore (1) holds with $M = Y d\lambda$ and $Y = \sum_{j=1}^n Y^j$. For such counting processes, we can prove the following result.

Theorem 5. *Assume that there exist a positive integer k and a positive number κ' , both known, such that $N^*(\mathcal{X}) \leq k$ a.s., (29) holds and $\text{Var}\left[\int_I s Y^* d\lambda\right] \leq \kappa' \mathbb{E}\left[\int_I s Y^* d\lambda\right]$ for all intervals $I \subset \mathcal{X}$. Assume moreover that $\int_{\mathcal{X}} s d\lambda < +\infty$ and the aggregated process N satisfies (2). Let us choose a family \mathcal{M} satisfying Assumption **H** and weights $\{\Delta'_m, m \in \mathcal{M}\}$ such that*

$$(30) \quad \sum_{m \in \mathcal{M}} \exp[-\eta \Delta'_m] = \Sigma'(\eta) < +\infty \quad \text{for } \eta = k \left(k + \int_{\mathcal{X}} s d\lambda\right)^{-1}.$$

Then the estimator $\hat{s}_{\hat{m}}$ defined in Section 2.3 with $\text{pen}(m) \geq 16\delta|m|(k+\kappa') + 404k\Delta'_m$ satisfies

$$\begin{aligned} \mathbb{E} [H^2(\hat{s}_{\hat{m}}, s)] &\leq 390 \left(\mathbb{E} \left[\inf_{m \in \mathcal{M}} (H^2(s, S_m) + \text{pen}(m)) \right] + 404k\eta^{-1}[\Sigma'(\eta)]^2 \right) + \varepsilon \\ &\leq 390 \left(\inf_{m \in \mathcal{M}} \{ \mathbb{E} [H^2(s, S_m)] + \text{pen}(m) \} + 404k\eta^{-1}[\Sigma'(\eta)]^2 \right) + \varepsilon. \end{aligned}$$

In the last bound, $\mathbb{E} [H^2(s, S_m)]$ plays the role of a bias term which can be bounded in the following way. Let us set

$$S'_m = \left\{ t = \sum_{I \in m \cap \mathcal{J}} t_I \mathbb{1}_I \quad \text{with } t_I \geq 0 \text{ for all } I \in m \right\} \cap \mathcal{L},$$

where the t_I are now deterministic. Then $S'_m \subset S_m$, hence $H^2(s, S_m) \leq H^2(s, S'_m)$ and, for $t \in S'_m$,

$$H^2(s, t) = \int_{\mathcal{X}} (\sqrt{s} - \sqrt{t})^2 Y d\lambda \leq n \int_{\mathcal{X}} (\sqrt{s} - \sqrt{t})^2 d\lambda,$$

since $Y \leq n$. Finally

$$\mathbb{E} [H^2(s, S_m)] \leq n \inf_{t \in S'_m} \int_{\mathcal{X}} (\sqrt{s} - \sqrt{t})^2 d\lambda = b_m^2(s)$$

and

$$\mathbb{E} [H^2(\hat{s}_{\hat{m}}, s)] \leq 390 \left(\inf_{m \in \mathcal{M}} \{ b_m^2(s) + \text{pen}(m) \} + \frac{404k}{\eta} [\Sigma'(\eta)]^2 \right) + \varepsilon.$$

Note that the present framework includes, as a particular case, density estimation, if we observe an n -sample X_1, \dots, X_n with density s with respect to λ and set $N^j(A) = \mathbb{1}_A(X_j)$. Then $Y = n$ and $H^2(s, t) = n \int_{\mathcal{X}} (\sqrt{s} - \sqrt{t})^2 d\lambda$ which corresponds to using the distance H of Section 4.1 multiplied by \sqrt{n} . Up to this scaling factor, the previous risk bound is analogue to that for estimating densities we get in Theorem 1.

In order to derive risk bounds which are similar to those given in Proposition 1, we have to distinguish between two situations. The most favorable one occurs when we know an upper bound Γ for $\int_{\mathcal{X}} s d\lambda$, in which case, since $0 \leq Y^* \leq 1$,

$$\text{Var} \left[\int_I s Y^* d\lambda \right] \leq \mathbb{E} \left[\left(\int_I s Y^* d\lambda \right)^2 \right] \leq \left(\int_{\mathcal{X}} s d\lambda \right) \mathbb{E} \left[\int_I s Y^* d\lambda \right]$$

and we can set $\kappa' = \Gamma$. Moreover, assuming that (6) holds, we can choose $\Delta'_m = (1 + k^{-1}\Gamma) \Delta_m$ without any further restriction on the family of models. Using the same family of partitions as in the density case, we recover the bounds of Propositions 1 and 2 up to the factor n corresponding to the rescaling of the distance H .

Let us now turn to the less favorable situation where no bound for $\int_{\mathcal{X}} s d\lambda$ is known, which is the typical case for Problem 4. As we shall see the number κ' can still be computed. As to (30) it will be satisfied with $\Delta'_m = |m|$ as soon as the number of models such that $|m| = D$ is bounded independently of D . Restricting ourselves to the family \mathcal{M}_R of regular partitions, we recover, up to the factor n , the bounds provided by case ii) of Proposition 1.

5.1. Survival analysis with right-censored data. Let us now consider the framework of Problem 3, denoting by P_T the common distribution of the T_i . We consider the counting process N on \mathbb{R}_+ defined by $N = \sum_{j=1}^n N^j$ where $N^j(A) = \mathbb{1}_{\{\tilde{T}_j \in A, D_j=1\}}$ for all measurable subsets A of \mathbb{R}_+ , so that we can take $k = 1$. Then the variables $N^j(A)$, $1 \leq j \leq n$ are i.i.d. Bernoulli random variables. We define s to be the hazard rate of the survival times, i.e. $s(t) = p(t)/\mathbb{P}[T_1 \geq t]$ for $t > 0$. Since s is not integrable on \mathbb{R}_+ we shall restrict ourselves to some bounded interval \mathcal{X} of \mathbb{R}_+ , which we can take, without loss of generality, to be $[0, 1)$ if we assume that $\mathbb{P}[T_1 \geq 1] > 0$. We also assume here that the censorship satisfies for all $t \geq 0$,

$$(31) \quad \mathbb{E} [N^j([0, t])] = \mathbb{E} \left[\int_0^t s(u) Y^j(u) du \right], \quad \text{with } Y^j(t) = \mathbb{1}_{\tilde{T}_j \geq t},$$

which means that (29) holds. Equality (31) is clearly satisfied when $C_j = T_j$ for all j , i.e. when the data are uncensored. It is also satisfied when the censorship is independent of the survival time, i.e. when C_j and T_j are independent for all j . Indeed, we then have for all j and $t \geq 0$, by Fubini Theorem and independence,

$$\begin{aligned} \mathbb{E} \left[\int_0^t s(u) Y^j(u) du \right] &= \mathbb{E} \left[\int_0^t \frac{p(u)}{\mathbb{P}(T_j \geq u)} \mathbb{1}_{C_j \geq u} \mathbb{1}_{T_j \geq u} du \right] \\ &= \int_0^t \frac{p(u) \mathbb{P}(T_j \geq u) \mathbb{P}(C_j \geq u)}{\mathbb{P}(T_j \geq u)} du \\ &= \int \mathbb{1}_{[0, t]}(u) \mathbb{P}(C_j \geq u) dP_T(u) \\ &= \mathbb{P}[T_j \leq t, T_j \leq C_j] = \mathbb{E} [N^j([0, t])] . \end{aligned}$$

Proposition 8. *If the processes N^j satisfy (31), the assumptions of Theorem 5 hold with $k = 1$, $\kappa' = 2$ and $\int_{\mathcal{X}} s d\lambda = -\log(\mathbb{P}[T_1 \geq 1])$.*

From a practical point of view, one can always estimate $\mathbb{P}[T_1 \geq 1]$ accurately enough to assume that an upper bound Γ for $\int_{\mathcal{X}} s d\lambda$ is known. We can therefore apply Theorem 5 to the family of models of Example 1 with the weights Δ_m given in Section 4.1, setting $\Delta'_m = (1 + \Gamma)\Delta_m$. We then obtain perfect analogues of Propositions 1 and 2 with constants C now depending on Γ . To avoid redundancy, we leave the precise statement of the risk bounds to the reader.

5.2. Transition intensities of Markov processes. Within the framework of Problem 4, we associate to $T_{1,0}$ the counting process N^* defined for $t \geq 0$ by $N^*([0, t]) = \mathbb{1}_{\{T_{1,0} \leq t\}}$ so that

$$(32) \quad \mathbb{E} [N^*([0, t])] = \int_0^t p(u) du = \mathbb{E} \left[\int_0^t \mathbb{1}_{\{X_{u-}=1\}} s(u) du \right]$$

and (29) holds with $Y^*(u) = \mathbb{1}_{\{X_{u-}=1\}}$. Our aim here is to estimate s on some bounded interval \mathcal{X} of \mathbb{R}_+ from the observation of the counting process $N = \sum_{j=1}^n N^j$ where the N^j 's are i.i.d. copies of N^* associated to n i.i.d. copies X^1, \dots, X^n of the process X . If X takes only the two values 0 and 1 and a.s. starts from 1 to reach 0, then the problem reduces to estimating the density p of $T_{1,0}$; it becomes novel when we have at least three states. In any case, we get the following result.

Proposition 9. *If the weights Δ'_m satisfy $\sum_{m \in \mathcal{M}} \exp[-\eta \Delta'_m] < +\infty$ for all $\eta > 0$ and $\int_{\mathcal{X}} s(t) dt < +\infty$ then Theorem 5 applies with $k = 1$ and $\kappa' = 2$.*

6. A UNIFYING RESULT

We want here to analyze our estimation procedure from the general point of view described in Section 2 and prove a risk bound for the estimator \tilde{s} , from which we shall be able to derive the previous risk bounds corresponding to all the specific frameworks that we considered. For this we introduce the following approximation for s in S_m :

$$(33) \quad \bar{s}_m = \sum_{I \in m \cap \mathcal{J}} \frac{s_I}{\lambda(I)} \mathbb{1}_I \quad \text{with } s_I = \int_I s d\lambda.$$

We need here a bound for $H^2(\hat{s}_m, \bar{s}_m)$ which holds uniformly for $m \in \overline{\mathcal{M}}$. It takes the following form:

H': There exist three positive constants a, b and $c, c \geq 1$ such that, for any $m \in \overline{\mathcal{M}}$,

$$(34) \quad \mathbb{P} [H^2(\hat{s}_m, \bar{s}_m) \geq c|m| + bz] \leq a \exp[-z] \quad \text{for all } z \geq 0.$$

We can now derive bounds for the risk of the estimator \tilde{s} defined in Section 2.3.

Theorem 6. *Let Assumptions **H** and **H'** hold and the weights Δ_m satisfy (6). Let the penalty $\text{pen}(m)$ be given by*

$$(35) \quad \text{pen}(m) \geq c\delta|m| + b\Delta_m.$$

and \hat{m} be any element of \mathcal{M} satisfying (5). Then the estimator $\tilde{s} = \hat{s}_{\hat{m}}$ satisfies

$$(36) \quad \mathbb{E} [H^2(\tilde{s}, s)] \leq 390 \left(\mathbb{E} \left[\inf_{m \in \mathcal{M}} (H^2(s, S_m) + \text{pen}(m)) \right] + ab\Sigma^2/2 \right) + \varepsilon.$$

Note that such a result has been obtained without any assumption on the underlying space \mathcal{X} and the true value s of the parameter, apart from the fact that it belongs to \mathcal{L} . Note also that in (36), the infimum over $m \in \mathcal{M}$ occurs inside the expectation, which makes a difference when M , and therefore $H(s, S_m)$, is random.

As we have previously seen, $\delta \leq 2$ for all the models we consider. Moreover, we shall see in Sections 7.3.1, 7.4.1 and 7.5.1 that for Problems 0, 1 and 2, $a = 1$ and b and c take the form $b = b'C_P$ and $c = c'C_P$ where b' and c' are numerical constants and C_P depends of the problem we consider (for instance $C_P = n^{-1}$ for density estimation). If we choose $\text{pen}(m) = c_0 C_P(|m| + \Delta_m)$ for some suitable numerical constant c_0 and $\varepsilon \leq C_P$, it follows that (36) becomes

$$\begin{aligned} & \mathbb{E} [H^2(\tilde{s}, s)] \\ & \leq 390 \left(\mathbb{E} \left[\inf_{m \in \mathcal{M}} (H^2(s, S_m) + c_0 C_P(|m| + \Delta_m)) \right] + 2b'C_P \Sigma^2/2 \right) + C_P, \end{aligned}$$

which gives (7). If there is only one model m in the family \mathcal{M} , we can fix $\Delta_m = 0$, hence $\Sigma = 1$, which leads to (3).

Proof. Let m^* be an arbitrary element of \mathcal{M} . It follows from the definition of \mathcal{D} that for any $m \in \mathcal{M}$, $H^2(\hat{s}_m, \hat{s}_{m^*}) \leq \mathcal{D}(m) \vee \mathcal{D}(m^*)$. Therefore,

$$(37) \quad H^2(\hat{s}_{\hat{m}}, \hat{s}_{m^*}) \leq \mathcal{D}(\hat{m}) \vee \mathcal{D}(m^*) \leq \mathcal{D}(m^*) + \varepsilon/3,$$

by (5). It also follows from (4) that, if $T_{m, m^*} \leq 0$, then

$$(38) \quad H^2(\hat{s}_m, \hat{s}_{m \vee m^*}) - H^2(\hat{s}_{m^*}, \hat{s}_{m \vee m^*}) \leq 16[\text{pen}(m^*) - \text{pen}(m)].$$

Moreover

$$\begin{aligned}
H^2(\hat{s}_m, \hat{s}_{m \vee m^*}) - H^2(\hat{s}_{m^*}, \hat{s}_{m \vee m^*}) \\
&= \int \hat{s}_m d\lambda - \int \hat{s}_{m^*} d\lambda + 2 \int \left(\sqrt{\hat{s}_{m^*}} - \sqrt{\hat{s}_m} \right) \sqrt{\hat{s}_{m \vee m^*}} d\lambda \\
&= H^2(\hat{s}_m, \hat{s}_{m^*}) + 2 \int \left(\sqrt{\hat{s}_{m^*}} - \sqrt{\hat{s}_m} \right) \left(\sqrt{\hat{s}_{m \vee m^*}} - \sqrt{\hat{s}_{m^*}} \right) d\lambda,
\end{aligned}$$

hence, by (38) and Cauchy-Schwarz Inequality,

$$\begin{aligned}
H^2(\hat{s}_m, \hat{s}_{m^*}) \\
&\leq 16[\text{pen}(m^*) - \text{pen}(m)] + 2 \int \left(\sqrt{\hat{s}_m} - \sqrt{\hat{s}_{m^*}} \right) \left(\sqrt{\hat{s}_{m \vee m^*}} - \sqrt{\hat{s}_{m^*}} \right) d\lambda \\
&\leq 16[\text{pen}(m^*) - \text{pen}(m)] + 2H(\hat{s}_m, \hat{s}_{m^*})H(\hat{s}_{m \vee m^*}, \hat{s}_{m^*}) \\
&\leq 16[\text{pen}(m^*) - \text{pen}(m)] + \frac{1}{2}H^2(\hat{s}_m, \hat{s}_{m^*}) + 4H^2(\hat{s}_{m \vee m^*}, \hat{s}_{m^*}).
\end{aligned}$$

Therefore, for any $m \in \mathcal{M}$ such that $T_{m, m^*} \leq 0$,

$$H^2(\hat{s}_m, \hat{s}_{m^*}) \leq 8H^2(\hat{s}_{m \vee m^*}, \hat{s}_{m^*}) + 32[\text{pen}(m^*) - \text{pen}(m)]$$

and, since

$$\begin{aligned}
H^2(\hat{s}_{m \vee m^*}, \hat{s}_{m^*}) \\
\leq 4 \left[H^2(\hat{s}_{m \vee m^*}, \bar{s}_{m \vee m^*}) + H^2(\bar{s}_{m \vee m^*}, s) + H^2(s, \bar{s}_{m^*}) + H^2(\bar{s}_{m^*}, \hat{s}_{m^*}) \right],
\end{aligned}$$

then

$$\begin{aligned}
(1/32)H^2(\hat{s}_m, \hat{s}_{m^*}) &\leq H^2(\hat{s}_{m \vee m^*}, \bar{s}_{m \vee m^*}) + H^2(\hat{s}_{m^*}, \bar{s}_{m^*}) + \text{pen}(m^*) \\
(39) \quad &\quad - \text{pen}(m) + H^2(\bar{s}_{m \vee m^*}, s) + H^2(s, \bar{s}_{m^*}).
\end{aligned}$$

Let us set, for all $z \geq 0$ and $(m, m') \in \mathcal{M}^2$,

$$\Omega_z = \bigcap_{(m, m') \in \mathcal{M}^2} \left\{ \omega \in \Omega \mid H^2(\hat{s}_{m \vee m'}, \bar{s}_{m \vee m'}) \leq c|m \vee m'| + b[\Delta_m + \Delta_{m'} + z] \right\}.$$

It follows from (34) that

$$(40) \quad \mathbb{P}[\Omega_z^c] \leq ae^{-z} \sum_{(m, m') \in \mathcal{M}^2} e^{-\Delta_m - \Delta_{m'}} = \Sigma^2 ae^{-z}.$$

Let now ω belong to Ω_z . It then follows that

$$(41) \quad H^2(\hat{s}_{m^*}, \bar{s}_{m^*}) \leq c|m^*| + 2b\Delta_{m^*} + bz \leq 2\text{pen}(m^*) + bz$$

and, using Assumption **H**, that

$$H^2(\hat{s}_{m \vee m^*}, \bar{s}_{m \vee m^*}) \leq c\delta[|m| + |m^*|] + b[\Delta_m + \Delta_{m^*} + z].$$

Therefore we derive from (39), (41) and (35) that, for all $m \in \mathcal{M}$ such that $T_{m, m^*} \leq 0$,

$$\begin{aligned}
(1/32)H^2(\hat{s}_m, \hat{s}_{m^*}) &\leq H^2(\bar{s}_{m \vee m^*}, s) + H^2(s, \bar{s}_{m^*}) + (1 + \delta)c|m^*| \\
&\quad + 3b\Delta_{m^*} + 2bz + \text{pen}(m^*) \\
&\leq H^2(\bar{s}_{m \vee m^*}, s) + H^2(s, \bar{s}_{m^*}) + 2bz + 4\text{pen}(m^*).
\end{aligned}$$

In order to control the bias terms $H^2(s, \bar{s}_{m'})$ of the various estimators involved in the construction of \tilde{s} , we shall use Lemma 2 below. Since $S_{m \vee m^*} \supset S_{m^*}$ for all $m \in \overline{\mathcal{M}}$, this lemma implies that

$$H^2(\bar{s}_{m' \vee m^*}, s) \leq 2H^2(s, S_{m' \vee m^*}) \leq 2H^2(s, S_{m^*}),$$

therefore

$$H^2(\hat{s}_m, \hat{s}_{m^*}) \leq 128 [H^2(s, S_{m^*}) + \text{pen}(m^*) + bz/2],$$

for all $m \in \mathcal{M}$ such that $T_{m, m^*} \leq 0$ and we conclude from (37) and the definition of \mathcal{D} that, if $\omega \in \Omega_z$,

$$H^2(\hat{s}_{\hat{m}}, \hat{s}_{m^*}) \leq \mathcal{D}(m^*) + \varepsilon/3 \leq 128 [H^2(s, S_{m^*}) + \text{pen}(m^*) + bz/2] + \varepsilon/3.$$

Since

$$H^2(\hat{s}_{\hat{m}}, s) \leq 3 [H^2(\hat{s}_{\hat{m}}, \hat{s}_{m^*}) + H^2(\hat{s}_{m^*}, \bar{s}_{m^*}) + H^2(\bar{s}_{m^*}, s)],$$

it follows from (41) and Lemma 2 that

$$H^2(\hat{s}_{\hat{m}}, s) \leq 3 [130H^2(s, S_{m^*}) + 130 \text{pen}(m^*) + 65bz + \varepsilon/3].$$

Since m^* is arbitrary in \mathcal{M} we finally get

$$H^2(\hat{s}_{\hat{m}}, s) \mathbb{1}_{\Omega_z} \leq 390 \left(\inf_{m \in \mathcal{M}} [H^2(s, S_m) + \text{pen}(m)] + bz/2 \right) + \varepsilon.$$

An integration with respect to z taking (40) into account leads to (36). \square

Lemma 2. *Within the framework of Section 2.1, for any $f \in \mathcal{L}$, we have*

$$H^2(f, \bar{f}_m) \leq 2H^2(f, S_m) \quad \text{with } \bar{f}_m = \sum_{I \in m \cap \mathcal{J}} \left(\int_I f \frac{d\lambda}{\lambda(I)} \right) \mathbb{1}_I.$$

Proof. Let $\mathcal{X}' = \bigcup_{I \in m \cap \mathcal{J}} I$. Note that M is a finite measure on \mathcal{X}' and that for all $t \in S_m$,

$$H^2(f, t) = H^2(f \mathbb{1}_{\mathcal{X}'}, t) + \int_{\mathcal{X} \setminus \mathcal{X}'} f d\lambda.$$

It is therefore enough to show the result for \mathcal{X}' in place of \mathcal{X} and $f \mathbb{1}_{\mathcal{X}'}$ in place of f and we can restrict ourselves to the case where M is a finite measure on \mathcal{X} . Let $\sqrt{f'}$ be the $\mathbb{L}_2(\mathcal{X}, d\lambda)$ projection of \sqrt{f} on S_m . Since the value of $\sqrt{f'}$ on I is given by $\int_I \sqrt{f} d\lambda / \lambda(I)$, it suffices to prove that for each $I \in m \cap \mathcal{J}$

$$(42) \quad \int_I \left(\sqrt{f} - \sqrt{\int_I f \frac{d\lambda}{\lambda(I)}} \right)^2 d\lambda \leq 2 \int_I \left(\sqrt{f} - \int_I \sqrt{f} \frac{d\lambda}{\lambda(I)} \right)^2 d\lambda.$$

By homogeneity, we may assume that $\lambda(I) = 1$. Expanding the left-hand side of (42) we get

$$\int_I \left(\sqrt{f} - \sqrt{\int_I f d\lambda} \right)^2 d\lambda = 2 \left(\int_I f d\lambda - \int_I \sqrt{f} d\lambda \times \sqrt{\int_I f d\lambda} \right),$$

which, together with the inequality $\sqrt{\int_I f d\lambda} \geq \int_I \sqrt{f} d\lambda$, leads to the desired result. \square

7. PROOFS

7.1. Proof of Lemma 1. Let $m = m_p \vee \mathcal{K}_j$ and $m' = m_{p'} \vee \mathcal{K}_{j'}$ be two elements of \mathcal{M} and $\bar{I}_p = ([0, 1]^k \setminus \cup_{I \in p} I)$, $\bar{I}_{p'} = ([0, 1]^k \setminus \cup_{I' \in p'} I')$. Assuming, with no loss of generality, that $j \geq j'$, we get

$$m \vee m' = m_p \vee m_{p'} \vee \mathcal{K}_j \vee \mathcal{K}_{j'} = m_p \vee m_{p'} \vee \mathcal{K}_j = m_1 \cup m_2 \cup m_3 \cup m_4,$$

with

$$\begin{aligned} m_1 &= \{K \cap I \cap I' \neq \emptyset \mid K \in \mathcal{K}_j, I \in p, I' \in p'\}; \\ m_2 &= \{K \cap I \cap \bar{I}_{p'} \neq \emptyset \mid K \in \mathcal{K}_j, I \in p\}; \\ m_3 &= \{K \cap \bar{I}_p \cap I' \neq \emptyset \mid K \in \mathcal{K}_j, I' \in p'\}; \\ m_4 &= \{K \cap \bar{I}_p \cap \bar{I}_{p'} \neq \emptyset \mid K \in \mathcal{K}_j\}. \end{aligned}$$

Since $j < J(p)$, hence $p \subset \cup_{l > j} \mathcal{K}_l$, for $K \in \mathcal{K}_j$ and $I \in p$, $K \cap I$ is either I or \emptyset , so that $m = p \cup p_j$ with $p_j = \{K \cap \bar{I}_p \neq \emptyset, K \in \mathcal{K}_j\}$ and $|m| = |p| + |p_j|$. It also follows that $|m_1| \leq |p| + |p'|$ and $|m_2| \leq |p|$. Then, given $K \in \mathcal{K}_j$ and $I' \in p'$, $K \cap I'$ is either K or I' or \emptyset since $K, I' \in \mathcal{K}$, so that $|m_3| \leq |p_j| + |p'|$. Finally $|m_4| \leq |p_j|$ and

$$|m \vee m'| \leq 2(|p| + |p'| + |p_j|) \leq 2(|m| + |m'|).$$

7.2. Some large deviations inequalities. The proofs of Theorems 1, 3, 4 and 5 require to check (34) for each specific framework. Since

$$(43) \quad H^2(\hat{s}_m, \bar{s}_m) = \sum_{I \in m \cap \mathcal{J}} \left(\sqrt{N(I)} - \sqrt{s_I} \right)^2 \quad \text{for all } m \in \mathcal{M},$$

this amounts to proving some deviation results for quantities of the form

$$\sum_{I \in m \cap \mathcal{J}} \left(\sqrt{N(I)} - \sqrt{s_I} \right)^2 - c|m|$$

which is the purpose of this section. Throughout it, we consider a finite set of non-negative random variables X_I with $I \in m$ and the related quantities

$$(44) \quad \chi^2(m) = \sum_{I \in m} \left(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]} \right)^2,$$

the notation suggesting that these variables behave roughly like χ^2 random variables as we shall see. Our purpose will be to derive deviation bounds for those variables from their expectation. Our first result is as follows:

Theorem 7. *Let $(X_I)_{I \in m}$ be a finite set of independent non-negative random variables and $\chi^2(m)$ be given by (44). We assume that there exists $\kappa > 0$ and $\tau \geq 0$ such that*

$$(45) \quad \log \left(\mathbb{E} \left[e^{z(X_I - \mathbb{E}[X_I])} \right] \right) \leq \kappa \frac{z^2 \mathbb{E}[X_I]}{2(1 - z\tau)} \quad \text{for all } z \in [0, 1/\tau[,$$

and

$$(46) \quad \log \left(\mathbb{E} \left[e^{-z(X_I - \mathbb{E}[X_I])} \right] \right) \leq \kappa \frac{z^2 \mathbb{E}[X_I]}{2} \quad \text{for all } z > 0.$$

Let

$$K = \max \left\{ \sqrt{2} ; \frac{\sqrt{2}}{2} + \sqrt{\left(\frac{\tau}{\kappa} - \frac{1}{2}\right)_+} \right\}.$$

Then for all $x > 0$,

$$(47) \quad \mathbb{P} \left[\chi^2(m) \geq \mathbb{E} [\chi^2(m)] + K^2 \kappa \left(2\sqrt{2|m|x} + x \right) \right] \leq e^{-x},$$

and

$$(48) \quad \mathbb{P} \left[\chi^2(m) \leq \mathbb{E} [\chi^2(m)] - 2K^2 \kappa \sqrt{2|m|x} \right] \leq e^{-x}.$$

Proof. Let us first introduce the following large deviation result, the proof of which follows the lines of the proof of Lemma 8 of Birgé and Massart (1998).

Lemma 3. *Let Y_1, \dots, Y_n be n independent, centered random variables. If*

$$\log (\mathbb{E} [e^{zY_i}]) \leq \kappa \frac{z^2 \theta_i}{2(1-z\tau)} \quad \text{for all } z \in [0, 1/\tau[\quad \text{and} \quad 1 \leq i \leq n,$$

then

$$\mathbb{P} \left[\sum_{i=1}^n Y_i \geq \left(2\kappa x \sum_{i=1}^n \theta_i \right)^{1/2} + \tau x \right] \leq e^{-x} \quad \text{for all } x > 0.$$

If, for $1 \leq i \leq n$ and all $z > 0$, $\log (\mathbb{E} [e^{-zY_i}]) \leq \kappa z^2 \theta_i / 2$, then

$$\mathbb{P} \left[\sum_{i=1}^n Y_i \leq - \left(2\kappa x \sum_{i=1}^n \theta_i \right)^{1/2} \right] \leq e^{-x} \quad \text{for all } x > 0.$$

It follows from (45), (46) and Lemma 3 with $n = 1$, $Y_1 = X_I - \mathbb{E} [X_I]$ and $\theta_1 = \mathbb{E} [X_I]$ that, for all $x > 0$ and $I \in m$,

$$\mathbb{P} \left[X_I \geq \mathbb{E} [X_I] + \sqrt{2\kappa \mathbb{E} [X_I] x} + \tau x \right] \leq e^{-x}$$

and

$$\mathbb{P} \left[X_I \leq \mathbb{E} [X_I] - \sqrt{2\kappa \mathbb{E} [X_I] x} \right] \leq e^{-x}.$$

Setting $u = \mathbb{E} [X_I] / (\kappa x)$, we deduce that, with probability not smaller than $1 - 2e^{-x}$,

$$\begin{aligned} & \left| \sqrt{X_I} - \sqrt{\mathbb{E} [X_I]} \right| \\ & \leq \max \left\{ \sqrt{\mathbb{E} [X_I]} - \sqrt{\left(\mathbb{E} [X_I] - \sqrt{2\kappa \mathbb{E} [X_I] x} \right)_+}; \right. \\ & \quad \left. \sqrt{\mathbb{E} [X_I] + \sqrt{2\kappa \mathbb{E} [X_I] x} + \tau x} - \sqrt{\mathbb{E} [X_I]} \right\} \\ & = \sqrt{\kappa x} \max \left\{ \sqrt{u} - \sqrt{\left(u - \sqrt{2u} \right)_+}; \sqrt{u + \sqrt{2u} + (\tau/\kappa)} - \sqrt{u} \right\} \\ & \leq \sqrt{\kappa x} \sup_{z>0} \max \left\{ \sqrt{z} - \sqrt{\left(z - \sqrt{2z} \right)_+}; \sqrt{z + \sqrt{2z} + (\tau/\kappa)} - \sqrt{z} \right\}. \end{aligned}$$

On the one hand, note that $z \rightarrow \sqrt{z} - \sqrt{(z - \sqrt{2z})_+}$ admits a maximum equal to $\sqrt{2}$ for $z = 2$. On the other hand, using the inequality $\sqrt{a+b} \leq \sqrt{a} + \sqrt{b}$ which holds for all positive numbers a, b , we obtain for all $z > 0$,

$$\begin{aligned} \sqrt{z + \sqrt{2z} + (\tau/\kappa)} - \sqrt{z} &\leq \sqrt{\left(\sqrt{z} + \frac{\sqrt{2}}{2}\right)^2 + \left(\frac{\tau}{\kappa} - \frac{1}{2}\right)_+} - \sqrt{z} \\ &\leq \frac{\sqrt{2}}{2} + \sqrt{\left(\frac{\tau}{\kappa} - \frac{1}{2}\right)_+} \end{aligned}$$

and therefore $|\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]}| \leq K\sqrt{\kappa x}$ with probability not smaller than $1 - 2e^{-x}$, or equivalently

$$(49) \quad \mathbb{P}[U_I \geq K^2 x] \leq 2e^{-x} \quad \text{for all } x > 0 \quad \text{with} \quad U_I = \kappa^{-1} \left(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]} \right)^2.$$

Since $\chi^2(m) = \kappa \sum_{I \in m} U_I$ and the random variables $U_I, I \in m$ are independent, (47) will derive from Lemma 3 if we show, setting $E_I = \mathbb{E}[U_I]$, that

$$(50) \quad \log \left(\mathbb{E} \left[e^{z(U_I - E_I)} \right] \right) \leq \frac{4K^4 z^2}{2(1 - K^2 z)} \quad \text{for all } z \in]0, 1/K^2[.$$

Similarly, (48) will follow from

$$(51) \quad \log \left(\mathbb{E} \left[e^{-z(U_I - E_I)} \right] \right) \leq \frac{4K^4 z^2}{2} \quad \text{for all } z > 0.$$

To prove (50), we shall use the following lemma about the centered moments of positive random variables.

Lemma 4. *Let Z be a non-negative random variable. For any positive even integer k ,*

$$\mathbb{E} \left[(Z - \mathbb{E}[Z])^k \right] \leq \mathbb{E} \left[Z^k \right] - (\mathbb{E}[Z])^k \leq \mathbb{E} \left[Z^k \right].$$

Note that the inequality $\mathbb{E} \left[(Z - \mathbb{E}[Z])^k \right] \leq \mathbb{E} \left[Z^k \right]$ also holds true for odd integers k since $\mathbb{E}[Z] \geq 0$ and the map $z \mapsto z^k$ is then increasing.

Proof. Since the result is trivial for $k = 2$, we may assume that $k \geq 4$ and, using homogeneity, that $\mathbb{E}[Z] = 1$. Consider the function $z \mapsto Q(z) = z^k - (z-1)^k - k(z-1)$ on $[0, +\infty[$. Its second derivative is negative for $z < 1/2$ and positive for $z > 1/2$, from which we easily derive that Q has a minimum for $z = 1$. This shows that $Q(z) \geq 1$ for all $z \geq 0$ and consequently,

$$\mathbb{E} \left[Z^k \right] - \mathbb{E} \left[(Z - 1)^k \right] = \mathbb{E} [Q(Z)] \geq Q(1) = 1$$

which leads to the result. \square

The random variable U_I is positive and by (49) satisfies $\mathbb{P}[U_I \geq t] \leq 2e^{-t/K^2}$. Consequently, we deduce from the previous lemma (with $Z = U_I$) that for all integers k (odd or even)

$$(52) \quad \mathbb{E} \left[(U_I - E_I)^k \right] \leq \mathbb{E} \left[U_I^k \right] = \int_0^{+\infty} k t^{k-1} \mathbb{P}[U_I \geq t] dt \leq 2(k!) K^{2k}.$$

Hence, for all $z \in]0, 1/K^2[$,

$$\log \left(\mathbb{E} \left[e^{z(U_I - E_I)} \right] \right) \leq \log \left(1 + 0 + 2 \sum_{k \geq 2} z^k K^{2k} \right) \leq 2 \sum_{k \geq 2} z^k K^{2k} = \frac{4K^4 z^2}{2(1 - K^2 z)}.$$

To prove (51), note that, for all $z, u > 0$, $e^{-zu} \leq 1 - zu + z^2 u^2 / 2$. Therefore, by (52),

$$\log \left(\mathbb{E} \left[e^{-z(U_I - E_I)} \right] \right) = \log \left(\mathbb{E} \left[e^{-zU_I} \right] \right) + zE_I \leq \frac{z^2}{2} \mathbb{E} [U_I^2] \leq \frac{4K^4 z^2}{2},$$

which completes the proof of Theorem 7. \square

A second pair of deviation inequalities for variables of the form $\chi^2(m)$ is as follows.

Theorem 8. *Let m be a finite index set and $\mathbf{X}_j = (X_{I,j})_{I \in m}$, $1 \leq j \leq p$ be i.i.d. random vectors with values in $\mathbb{R}_+^{|m|}$. Assume that there exist positive numbers A and κ such that*

$$(53) \quad \sum_{I \in m} X_{I,1} \leq A \text{ a.s.} \quad \text{and} \quad \text{Var}(X_{I,1}) \leq \kappa \mathbb{E}[X_{I,1}] \quad \text{for all } I \in m.$$

If $X_I = \sum_{j=1}^p X_{I,j}$ for all $I \in m$ and $\chi^2(m)$ is given by (44), then

$$(54) \quad \mathbb{P} [\chi^2(m) \geq 8\kappa|m| + 202Ax] \leq e^{-x} \quad \text{for all } x > 0.$$

Proof. Since $X_{I,1} = 0$ a.s. if $\mathbb{E}[X_{I,1}] = 0$, we may remove all indexes I such that $\mathbb{E}[X_{I,1}] = 0$ in the sum and therefore assume that $\mathbb{E}[X_I] = p\mathbb{E}[X_{I,1}] > 0$ for all $I \in m$. We can then write, for all $z > 0$,

$$\begin{aligned} & \mathbb{P} \left(\sqrt{\chi^2(m)} \geq z \right) \\ &= \mathbb{P} \left(\sum_{I \in m} \frac{(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]})}{\sqrt{\chi^2(m)}} \left(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]} \right) \geq z, \sqrt{\chi^2(m)} \geq z \right) \\ &= \mathbb{P} \left(\sum_{I \in m} \frac{(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]})}{\sqrt{\chi^2(m)}} \frac{X_I - \mathbb{E}[X_I]}{\sqrt{X_I} + \sqrt{\mathbb{E}[X_I]}} \geq z, \sqrt{\chi^2(m)} \geq z \right) \\ &= \mathbb{P} \left(\sum_{j=1}^p \left[\sum_{I \in m} \frac{(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]}) (X_{I,j} - \mathbb{E}[X_{I,j}])}{\sqrt{\chi^2(m)} (\sqrt{X_I} + \sqrt{\mathbb{E}[X_I]})} \right] \geq z, \sqrt{\chi^2(m)} \geq z \right) \\ &= \mathbb{P} \left(\sum_{I \in m} \left[\sum_{j=1}^p \frac{(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]}) \sqrt{\mathbb{E}[X_I]} X_{I,j} - \mathbb{E}[X_{I,j}]}{\sqrt{\chi^2(m)} (\sqrt{X_I} + \sqrt{\mathbb{E}[X_I]}) \sqrt{\mathbb{E}[X_I]}} \right] \geq z, \sqrt{\chi^2(m)} \geq z \right) \\ &= \mathbb{P} \left(\sum_{j=1}^p \sum_{I \in m} t_I \frac{X_{I,j} - \mathbb{E}[X_{I,j}]}{\sqrt{\mathbb{E}[X_I]}} \geq z, \sqrt{\chi^2(m)} \geq z \right), \end{aligned}$$

where

$$t_I = \frac{(\sqrt{X_I} - \sqrt{\mathbb{E}[X_I]}) \sqrt{\mathbb{E}[X_I]}}{\sqrt{\chi^2(m)} (\sqrt{X_I} + \sqrt{\mathbb{E}[X_I]})} \quad \text{for all } I \in m.$$

Note that $\sum_{I \in m} t_I^2 \leq 1$ since $\sqrt{\mathbb{E}[X_I]} / (\sqrt{X_I} + \sqrt{\mathbb{E}[X_I]}) \leq 1$ and that $|t_I| \leq z^{-1} \sqrt{\mathbb{E}[X_I]}$ on the set $\sqrt{\mathcal{X}^2(m)} \geq z$, from which we deduce that

$$(55) \quad \mathbb{P} \left(\sqrt{\mathcal{X}^2(m)} \geq z \right) \leq \mathbb{P} \left(\sup_{\mathbf{t} \in \mathcal{T}} \sum_{j=1}^p \sum_{I \in m} t_I \frac{X_{I,j} - \mathbb{E}[X_{I,j}]}{\sqrt{\mathbb{E}[X_I]}} \geq z \right),$$

where \mathcal{T} denotes the set of vectors $\mathbf{t} = (t_I)_{I \in m} \in \mathbb{R}^{|m|}$ satisfying

$$(56) \quad |t_I| \leq \frac{\sqrt{\mathbb{E}[X_I]}}{z} \quad \text{for all } I \in m \quad \text{and} \quad \sum_{I \in m} t_I^2 \leq 1.$$

In order to bound the right-hand side of (55), we shall use the following result from Massart (2000, Theorem 2.4).

Theorem 9. *Let ξ_1, \dots, ξ_p be independent random variables with values in some measurable space \mathcal{H} and \mathcal{F} be some countable family of real valued measurable functions on \mathcal{H} such that $\|f\|_\infty \leq b < +\infty$ for all $f \in \mathcal{F}$. If*

$$Z = \sup_{f \in \mathcal{F}} \left| \sum_{j=1}^p f(\xi_j) - \mathbb{E}[f(\xi_j)] \right| \quad \text{and} \quad \sigma^2 = \sup_{f \in \mathcal{F}} \left[\sum_{j=1}^p \text{Var}(f(\xi_j)) \right],$$

then for every positive numbers ε, x

$$\mathbb{P} \left[Z \geq (1 + \varepsilon) \mathbb{E}[Z] + 2\sigma\sqrt{2x} + (2.5 + 32\varepsilon^{-1})bx \right] \leq e^{-x}.$$

We want to apply this result to the vectors $\xi_j \in \mathbb{R}^{|m|}$ with coordinates $\xi_{I,j} = (X_{I,j} - \mathbb{E}[X_{I,j}]) / \sqrt{\mathbb{E}[X_I]}$ for $I \in m$. Under our assumptions, these random vectors are independent and satisfy

$$\sum_{I \in m} \sqrt{\mathbb{E}[X_I]} |\xi_{I,j}| \leq \sum_{I \in m} (X_{I,j} + \mathbb{E}[X_{I,j}]) \leq 2A.$$

Consequently, the random vectors ξ_j take their values in the subset \mathcal{H} of $\mathbb{R}^{|m|}$ given by

$$\mathcal{H} = \left\{ \mathbf{u} = (u_I, I \in m) \mid \sum_{I \in m} \sqrt{\mathbb{E}[X_I]} |u_I| \leq 2A \right\}.$$

For $\mathbf{u} \in \mathcal{H}$ and $\mathbf{t} \in \mathcal{T}$, we set $f_{\mathbf{t}}(\mathbf{u}) = \sum_{I \in m} t_I u_I$ and $\mathcal{F} = \{f_{\mathbf{t}}, \mathbf{t} \in \mathcal{T}'\}$ where \mathcal{T}' denotes a countable and dense subset of \mathcal{T} . With no loss of generality we can assume that \mathcal{T}' is symmetric around 0 (if $\mathbf{t} \in \mathcal{T}'$ then $-\mathbf{t} \in \mathcal{T}'$) which implies that the absolute values can be removed in the definition of Z . Since, for all $\mathbf{t} \in \mathcal{T}$ and $1 \leq j \leq p$, $f_{\mathbf{t}}(\xi_j)$ is centered, we can finally write

$$Z = \sup_{\mathbf{t} \in \mathcal{T}} \sum_{j=1}^p \sum_{I \in m} t_I \frac{X_{I,j} - \mathbb{E}[X_{I,j}]}{\sqrt{\mathbb{E}[X_I]}} = \sup_{\mathbf{t} \in \mathcal{T}} \sum_{I \in m} t_I \left(\sum_{j=1}^p \frac{X_{I,j} - \mathbb{E}[X_{I,j}]}{\sqrt{\mathbb{E}[X_I]}} \right).$$

Using Cauchy-Schwarz Inequality and (56), we then derive that

$$\mathbb{E}^2[Z] \leq \mathbb{E}[Z^2] \leq \sum_{I \in m} \mathbb{E} \left[\left(\sum_{j=1}^p \frac{X_{I,j} - \mathbb{E}[X_{I,j}]}{\sqrt{\mathbb{E}[X_I]}} \right)^2 \right] = \sum_{I \in m} \sum_{j=1}^p \frac{\text{Var}(X_{I,j})}{\mathbb{E}[X_I]}.$$

Since $\text{Var}(X_{I,j}) \leq \kappa \mathbb{E}[X_{I,j}]$ and $\sum_{j=1}^p \mathbb{E}[X_{I,j}] = \mathbb{E}[X_I]$, we conclude that $\mathbb{E}[Z] \leq \sqrt{\kappa|m|}$. To bound $\|f_{\mathbf{t}}\|_\infty$, we use (56) which implies that, for all $\mathbf{u} \in \mathcal{H}$ and $\mathbf{t} \in \mathcal{T}$,

$$|f_{\mathbf{t}}(\mathbf{u})| = \left| \sum_{I \in m} t_I u_I \right| \leq \sum_{I \in m} |t_I| |u_I| \leq \sum_{I \in m} \frac{\sqrt{\mathbb{E}[X_I]} |u_I|}{z} \leq \frac{2A}{z}.$$

Finally, it follows from the equidistribution of the \mathbf{X}_j , Cauchy-Schwarz Inequality, (53) and (56) that, for all $\mathbf{t} \in \mathcal{T}$,

$$\begin{aligned} \sum_{j=1}^p \text{Var}(f_{\mathbf{t}}(\xi_j)) &= p \text{Var}(f_{\mathbf{t}}(\xi_1)) = p \mathbb{E} \left[\left(\sum_{I \in m} t_I \frac{X_{I,1} - \mathbb{E}[X_{I,1}]}{\sqrt{p \mathbb{E}[X_{I,1}]}} \right)^2 \right] \\ &\leq 2 \mathbb{E} \left[\left(\sum_{I \in m} t_I \frac{X_{I,1}}{\sqrt{\mathbb{E}[X_{I,1}]}} \right)^2 \right] + 2 \left(\sum_{I \in m} t_I \sqrt{\mathbb{E}[X_{I,1}]} \right)^2 \\ &\leq 2 \mathbb{E} \left[\left(\sum_{I \in m} X_{I,1} \right) \left(\sum_{I \in m} t_I^2 \frac{X_{I,1}}{\mathbb{E}[X_{I,1}]} \right) \right] + 2 \sum_{I \in m} t_I^2 \sum_{I \in m} \mathbb{E}[X_{I,1}] \\ &\leq 2A \left(\mathbb{E} \left[\sum_{I \in m} t_I^2 \frac{X_{I,1}}{\mathbb{E}[X_{I,1}]} \right] + \sum_{I \in m} t_I^2 \right) \leq 4A. \end{aligned}$$

In view of all these bounds, we may apply Theorem 9 with $\sigma^2 = 4A$, $b = 2A/z$ and $\varepsilon = 1$ and obtain that $\mathbb{P}[\sqrt{\chi^2(m)} \geq z] \leq e^{-x}$ as soon as $z \geq 2\sqrt{\kappa|m|} + 4\sqrt{2Ax} + 69Ax/z$. Solving this quadratic inequation and using $(a+b)^2 \leq 2(a^2+b^2)$, we can check that this inequality holds if $z^2 \geq 8\kappa|m| + 202Ax$, hence the result. \square

7.3. Density estimation.

7.3.1. Proof of Theorem 1. For two given classes $m, m' \in \mathcal{M}$, we apply Theorem 8 with $m'' = m \vee m'$ in place of m , $p = n$ and $X_{I,j} = \mathbb{1}_{Y_j \in I}$ for all $I \in m''$ and $j = 1, \dots, n$. Then $X_I = nN(I)$ and (53) is satisfied with $A = \kappa = 1$ since $X_{I,1}$ is a Bernoulli random variable and we derive from (43) that, for all $x > 0$, with probability not smaller than $1 - e^{-x}$,

$$H^2(\hat{s}_{m''}, \bar{s}_{m''}) = \sum_{I \in m''} \left(\sqrt{N(I)} - \sqrt{\mathbb{E}[N(I)]} \right)^2 = \frac{\chi^2(m'')}{n} \leq \frac{8|m''| + 202x}{n}.$$

Therefore (34) holds with $c = 8/n$, $a = 1$ and $b = 202/n$. We then conclude from Theorem 6 and the fact that $H^2(t, u)$ is always bounded by 2.

7.3.2. Proof of Proposition 1. By assumption, \sqrt{s} has a variation bounded by R and we may apply to it Corollary 1 of Barron, Birgé and Massart (1999) with $\alpha = 1$, $D = 2^j$ with $j \geq 2$ and $N = 2^{3j}$. It follows that one can find $m \in \mathcal{M}_{3j,D}$ such that $H^2(s, S_m) \leq (64/3)(R/D)^2$. Since $\text{pen}(m) \leq CjDn^{-1}$ for $m \in \mathcal{M}_{3j,D}$, we derive from Theorem 1 that

$$\mathbb{E}_s [H^2(\tilde{s}, s)] \leq C' \inf_{j \geq 2} \{R^2 2^{-2j} + j 2^j n^{-1}\}.$$

Then (11) follows if we define $j \geq 2$ by

$$4^{-j+1} \leq [nR^2 / \log(1 + nR^2)]^{-2/3} < 4^{-j+2},$$

which is always possible since $nR^2 > 0$, and distinguish between the cases $j = 2$ (which corresponds to $nR^2 \leq 26.519$) and $j > 2$.

When \sqrt{s} is continuous with modulus w , there exists an element $t \in S_{m_j}$ such that $\|\sqrt{s} - \sqrt{t}\|_\infty \leq w(2^{-j})$, hence $H(s, S_{m_j}) \leq w(2^{-j})$. Since $x_w > 0$, we can choose j such that $2^{-j} < x_w \leq 2^{-j+1}$. Recalling that $\text{pen}(m_j) \leq C2^j/n$, we deduce from Theorem 1 that

$$\mathbb{E}_s [H^2(\tilde{s}, s)] \leq C' [w^2(2^{-j}) + n^{-1}2^j] \leq C' [w^2(x_w) + 2(nx_w)^{-1}] \leq 3C'(nx_w)^{-1},$$

which proves (12). If \sqrt{s} belongs to \mathcal{H}_α^R with $R \geq n^{-1/2}$, then $x_w = (nR)^{-2/(2\alpha+1)}$ and the risk bound follows.

If s belongs to $\mathcal{S}^3(D, R)$, we can write $s = \sum_{k=1}^D s_k \mathbb{1}_{[x_{k-1}, x_k]}$ with $0 = x_0 < x_1 < \dots < x_D = 1$ and $\sup_{1 \leq k \leq D} s_k \leq R$. Fix l such that $2^l \geq nR > 2^{l-1}$. Then $2^l \geq 2D$ and for $0 \leq k \leq D$, set $x'_k = \sup\{x \in \mathcal{J}_l \mid x \leq x_k\}$ and $t = \sum_{k=1}^D s_k \mathbb{1}_{[x'_{k-1}, x'_k]}$ so that $t \in S_m$ with $m \in \mathcal{M}_{l, D'}$ with $D' \leq D$ since some intervals $[x'_{k-1}, x'_k]$ may be empty. Then

$$H^2(s, t) \leq R \sum_{k=1}^{D-1} (x_k - x'_k) < RD2^{-l}.$$

Recalling from (9) that $\text{pen}(m) \leq Cn^{-1}[D(l \log 2 + 2 - \log D) + 2 \log l]$ for $m \in \mathcal{M}_{l, D}$, we conclude from Theorem 1, (9) and our choice of l that

$$\begin{aligned} \mathbb{E}_s [H^2(\tilde{s}, s)] &\leq C' \left[RD2^{-l} + [D(l \log 2 + 2 - \log D) + 2 \log l] n^{-1} \right] \\ &\leq C'(D/n) \left[3 + \log 2 + \log(2^{l-1}/D) + 2D^{-1} \log l \right] \\ &\leq C'(D/n) \left[3 + \log 2 + \log(nR/D) + 2(D \log 2)^{-1} \log \log(2nR) \right] \end{aligned}$$

and (13) follows since $nR \geq 2D$.

7.4. Random vectors.

7.4.1. Proof of Theorem 4. For two given elements $m, m' \in \mathcal{M}$, we apply Theorem 7 with $m'' = m \vee m'$ in place of m and $X_I = N(I)$. We derive from the independence of the N_i that (45) and (46) hold. Therefore, for all $x > 0$, with probability not smaller than $1 - e^{-x}$,

$$\begin{aligned} H^2(\hat{s}_{m''}, \bar{s}_{m''}) &= \sum_{I \in m''} \left(\sqrt{N(I)} - \sqrt{\mathbb{E}[N(I)]} \right)^2 \\ &\leq \mathbb{E} \left[\sum_{I \in m''} \left(\sqrt{N(I)} - \sqrt{\mathbb{E}[N(I)]} \right)^2 \right] + K^2 \kappa \left(2\sqrt{2|m''|x} + x \right). \end{aligned}$$

It follows from (45) that $\text{Var}(N(I)) \leq \kappa \mathbb{E}[N(I)]$ (expand both side of (45) in a vicinity of 0) and therefore

$$\begin{aligned} \mathbb{E} \left[\sum_{I \in m''} \left(\sqrt{N(I)} - \sqrt{\mathbb{E}[N(I)]} \right)^2 \right] &= \sum_{I \in m''} \mathbb{E} \left[\left(\sqrt{N(I)} - \sqrt{\mathbb{E}[N(I)]} \right)^2 \right] \\ &\leq \sum_{I \in m''} \mathbb{E} \left[\frac{(N(I) - \mathbb{E}[N(I)])^2}{\mathbb{E}[N(I)]} \right] \leq \kappa |m''|. \end{aligned}$$

Using the inequality $2\sqrt{2|m''|x} \leq |m''| + 2x$ we conclude that, with probability not smaller than $1 - e^{-x}$,

$$(57) \quad H^2(\hat{s}_{m''}, \bar{s}_{m''}) \leq (1 + K^2) \kappa |m''| + 3K^2 \kappa x.$$

We derive that (34) is fulfilled with $c = (1 + K^2) \kappa$, $b = 3K^2 \kappa$, $a = 1$ and Theorem 4 follows from Theorem 6.

7.4.2. Proof of Proposition 7. Let us first note that, if $|m| = n$, then $H^2(s, S_m) = 0$, hence by (28), $\mathbb{E}[H^2(\tilde{s}, s)] \leq C(\kappa, K)n$ which proves the bound when $R > n/\sqrt{3}$. For the other cases, we deduce from Lemma 5 below that, for any $D \in \mathcal{X}$, one can find some $m \in \mathcal{M}$ such that $|m| \leq D$ and $H^2(s, S_m) \leq n(R/D)^2$. Setting $D = 1$, we get the result for the case $R^2 < n^{-1} \log n$. Finally, when $n^{-1} \log n \leq R^2 \leq n^2/3$ we fix $D = \inf \{j \in \mathbb{N} \mid j^3 \geq nR^2/\log(n/R)\}$. Since the function $R \mapsto R^2/\log(n/R)$ is increasing for $R < n/\sqrt{3}$, $1 \leq D \leq n$ and the corresponding risk bound follows.

Lemma 5. *Let f be a nondecreasing function from $\mathcal{X} = \{1, \dots, n\}$ to \mathbb{R} such that $\sqrt{f(n)} - \sqrt{f(1)} = R$. For $D \in \mathcal{X}$, one can find a partition (I_1, \dots, I_K) of \mathcal{X} into $K \leq D$ intervals and a function g from \mathcal{X} to \mathbb{R} of the form $g = \sum_{k=1}^K \beta_k \mathbb{1}_{I_k}$ such that*

$$\sum_{i=1}^n \left(\sqrt{f(i)} - \sqrt{g(i)} \right)^2 \leq nR^2 D^{-2}.$$

Proof: Let us set $j_0 = 1$ and define iteratively for $k \geq 1$, using the convention $\inf \emptyset = n$,

$$(58) \quad j_k = \inf \left\{ j \in \{j_{k-1} + 1, \dots, n\} \mid \sqrt{f(j)} - \sqrt{f(j_{k-1})} > R/D \right\}.$$

Let $K = \inf \{k \geq 1, j_k = n\}$, $I_K = \{j_{K-1}, \dots, n\}$ and for $k = 1, \dots, K-1$ (if $K \geq 2$), $I_k = \{j_{k-1}, \dots, j_k - 1\}$. This defines a partition of \mathcal{X} with K elements and it follows from (58) that

$$R = \sqrt{f(n)} - \sqrt{f(1)} \geq \sum_{k=1}^{K-1} \sqrt{f(j_k)} - \sqrt{f(j_{k-1})} > (K-1)R/D,$$

hence $K-1 < D$ and $K \leq D$. Let us now set $\beta_k = f(j_{k-1})$ for $1 \leq k \leq K$. Since $\sqrt{f(j_k - 1)} - \sqrt{f(j_{k-1})} \leq R/D$ we get for all $i \in I_k$, $0 \leq \sqrt{f(i)} - \sqrt{g(i)} \leq R/D$. Hence,

$$\sum_{i=1}^n \left(\sqrt{f(i)} - \sqrt{g(i)} \right)^2 = \sum_{k=1}^K \sum_{i \in I_k} \left(\sqrt{f(i)} - \sqrt{g(i)} \right)^2 \leq nR^2 D^{-2}. \quad \square$$

7.5. Poisson and other counting processes.

7.5.1. Poisson processes. The proof of Theorem 3 follows the same lines as the proof of Theorem 4. We apply Theorem 7 with $m'' = m \vee m'$ in place of m and $X_I = N(I)$. Since $\{N(I), I \in m''\}$ are independent Poisson random variables, the assumptions of the theorem are fulfilled with $\kappa = \tau = 1$. We then proceed as for Theorem 4 to get (57) with $K^2 = 2$ which provides the relevant values of c and b .

7.5.2. *Proof of Theorem 5.* Let us fix two classes $m, m' \in \mathcal{M}$. We first apply Theorem 8 with $m'' = m \vee m'$ in place of m , $p = n$ and $X_{I,j} = N^j(I)$ for all $I \in m''$ and $j = 1, \dots, n$. Then for all $I \in m''$, $N(I) = X_I$. Since $X_{I,j}$ is bounded by k , $\mathbb{E}[X_{I,j}^2] \leq k\mathbb{E}[X_{I,j}]$ and (53) holds with $A = \kappa = k$. This implies that, for all $x > 0$, with probability not smaller than $1 - e^{-x}$,

$$(59) \quad \sum_{I \in m''} \left(\sqrt{N(I)} - \sqrt{\mathbb{E}[N(I)]} \right)^2 \leq k(8|m''| + 202x).$$

Then we apply once again Theorem 8 with $m'' = m \vee m'$ in place of m , $p = n$ and $X_{I,j} = \int_I sY^j d\lambda$ for all $I \in m''$ and $j = 1, \dots, n$. Since Y^j is bounded by 1, the assumptions of Theorem 8 are fulfilled with $A = \int_{\mathcal{X}} s d\lambda$ and $\kappa = \kappa'$. Consequently, with probability not smaller than $1 - e^{-x}$,

$$(60) \quad \sum_{I \in m''} \left(\sqrt{\int_I sY d\lambda} - \sqrt{\mathbb{E}\left[\int_I s d\lambda\right]} \right)^2 \leq 8\kappa'|m''| + 202Ax.$$

Since $\mathbb{E}[\int_I s d\lambda] = \mathbb{E}[N(I)]$, we derive from (59) and (60) that, with probability not smaller than $1 - 2e^{-x}$,

$$\begin{aligned} & H^2(\hat{s}_{m''}, \bar{s}_{m''}) \\ & \leq \sum_{I \in m''} \left(\sqrt{N(I)} - \sqrt{\int_I sY d\lambda} \right)^2 \\ & \leq 2 \sum_{I \in m''} \left(\sqrt{N(I)} - \sqrt{\mathbb{E}[N(I)]} \right)^2 + 2 \sum_{I \in m''} \left(\sqrt{\int_I sY d\lambda} - \sqrt{\mathbb{E}\left[\int_I s d\lambda\right]} \right)^2 \\ & \leq 16|m''|(k + \kappa') + 404x(k + A). \end{aligned}$$

This means that (34) holds with $c = 16(k + \kappa')$, $a = 2$ and $b = 404(k + A)$. Therefore, if we set $\Delta_m = k(k + A)^{-1}\Delta'_m$ for all $m \in \mathcal{M}$, (6) holds with $\Sigma = \Sigma'(k/(k + A))$ and $\text{pen}(m) = 16\delta|m|(k + \kappa') + 404k\Delta'_m$. An application of Theorem 6 leads to the result.

7.5.3. *Proof of Proposition 8.* The following argument shows that (2) is satisfied: let A be some measurable subset of \mathcal{X} and B be the subset of A given by $B = \{t \in A \mid \lambda([0, t] \cap A) = 0\}$. Since, by definition, the sets $[0, t] \cap B$ with $t \in B$ are negligible, $\lambda(B) = 0$ (write B as an at most countable union of those sets). Consequently,

$$\begin{aligned} \mathbb{P}(N(A) > 0, M(A) = 0) & \leq \sum_{j=1}^n \mathbb{P}\left(N^j(A) = 1, \int_A \mathbb{1}_{\tilde{T}_j \geq t} dt = 0\right) \\ & \leq \sum_{j=1}^n \mathbb{P}\left(\tilde{T}_j = T_j, T_j \in A, \lambda(A \cap [0, \tilde{T}_j]) = 0\right) \\ & \leq \sum_{j=1}^n \mathbb{P}(T_j \in B) = 0 \end{aligned}$$

since the common distribution of the T_j is continuous. Moreover

$$\int_{\mathcal{X}} s d\lambda = \int_0^1 \frac{p(t)}{\mathbb{P}[T_1 \geq t]} dt = -\log(\mathbb{P}(T_1 \geq 1))$$

since $-p(t)$ is the derivative of $\mathbb{P}[T_1 \geq t]$. Finally we can take $\kappa' = 2$ since, whatever $I \subset \mathcal{X}$,

$$\begin{aligned} \text{Var} \left[\int_I s(t) Y_t^* dt \right] &\leq \mathbb{E} \left[\left(\int_I s(t) Y_t^* dt \right)^2 \right] = \mathbb{E} \left[\int_{I \times I} s(t) s(t') Y_t^* Y_{t'}^* dt dt' \right] \\ &= \int_{I \times I} s(t) s(t') \mathbb{E} [Y_t^* Y_{t'}^*] dt dt' \\ &= \int_{I \times I} s(t) s(t') \mathbb{P} [\tilde{T}_1 \geq \max \{t, t'\}] dt dt' \\ &= 2 \int_I s(t) \left(\int_I \mathbb{1}_{\{t' \geq t\}} s(t') \mathbb{P} [\tilde{T}_1 \geq t'] dt' \right) dt \\ &\leq 2 \int_I s(t) \mathbb{E} \left[\int_t^1 s(t') Y_{t'}^* dt' \right] dt \\ &= 2 \int_I s(t) \mathbb{E} [N^1([t, 1])] dt \\ &\leq 2 \int_I s(t) \mathbb{P} [\tilde{T}_1 \geq t] dt = 2 \mathbb{E} \left[\int_I s(t) Y_t^* dt \right]. \end{aligned}$$

7.5.4. Proof of Proposition 9. Clearly (29) holds true. We now prove that Condition (2) is also fulfilled. Let A be some measurable subset of \mathbb{R}_+ and for $l \geq 1$ let B_l be the subset of A defined by

$$B_l = \{t \in A \mid \lambda([t - l^{-1}, t] \cap A) = 0\}.$$

For each $l \geq 1$, note that the sets $[t - l^{-1}, t] \cap B_l \subset [t - l^{-1}, t] \cap A$ are negligible for $t \in B_l$ and hence so is B_l (write B_l as an at most countable union of those). Denoting, for $j = 1, \dots, n$, the time of the jump of X^j from state 1 to 0 by $T_{1,0}^j$, we have

$$\begin{aligned} \mathbb{P}(N(A) > 0, M(A) = 0) &\leq \sum_{j=1}^n \mathbb{P} \left(N^j(A) = 1, \int_A \mathbb{1}_{X_{t-}^j = 1} dt = 0 \right) \\ &\leq \sum_{j=1}^n \mathbb{P} \left(N^j(A) = 1, \exists \varepsilon > 0, \lambda([T_{1,0}^j - \varepsilon, T_{1,0}^j] \cap A) = 0 \right) \\ &\leq \sum_{j=1}^n \sum_{l \geq 1} \mathbb{P} \left(T_{1,0}^j \in A, \lambda([T_{1,0}^j - l^{-1}, T_{1,0}^j] \cap A) = 0 \right) \\ &\leq \sum_{j=1}^n \sum_{l \geq 1} \mathbb{P} \left(T_{1,0}^j \in B_l \right) = \sum_{j=1}^n \sum_{l \geq 1} \mathbb{E} [N^*(B_l)] = 0, \end{aligned}$$

by (32). We may clearly fix $k = 1$ and the choice of κ' is justified by the following argument. First note that whatever $I \subset \mathcal{X}$ and $t > 0$

$$\begin{aligned}
& \mathbb{P}(X_{t-}^1 = 1, T_{1,0}^1 \in I, T_{1,0}^1 \geq t) \\
&= \int_I \mathbb{1}_{\{u \geq t\}} \mathbb{P}(X_{t-}^1 = 1, u \leq T_{1,0}^1 \leq u + du) \\
&= \int_I \mathbb{1}_{\{u \geq t\}} \mathbb{P}(X_{t-}^1 = 1, X_{u-}^1 = 1) \mathbb{P}(u \leq T_{1,0}^1 \leq u + du \mid X_{t-}^1 = 1, X_{u-}^1 = 1) \\
&= \int_I \mathbb{1}_{\{u \geq t\}} \mathbb{P}(X_{t-}^1 = 1, X_{u-}^1 = 1) \mathbb{P}(u \leq T_{1,0}^1 \leq u + du \mid X_{u-}^1 = 1)
\end{aligned}$$

since X^1 is a Markov process. Hence

$$\begin{aligned}
\mathbb{P}(X_{t-}^1 = 1, T_{1,0}^1 \in I, T_{1,0}^1 \geq t) &= \int_I \mathbb{1}_{\{u \geq t\}} \mathbb{P}(X_{t-}^1 = 1, X_{u-}^1 = 1) s(u) du \\
&= \mathbb{E} \left[\int_I \mathbb{1}_{\{u \geq t\}} \mathbb{1}_{\{X_{t-}^1 = 1\}} \mathbb{1}_{\{X_{u-}^1 = 1\}} s(u) du \right].
\end{aligned}$$

It then follows that

$$\begin{aligned}
\text{Var} \left(\int_I s(t) Y_t^1 dt \right) &\leq \mathbb{E} \left[\left(\int_{\mathcal{X}} \mathbb{1}_I(t) s(t) Y_t^1 dt \right)^2 \right] \\
&= \mathbb{E} \left[\int_{\mathcal{X} \times \mathcal{X}} \mathbb{1}_I(t) \mathbb{1}_I(u) s(t) s(u) Y_t^1 Y_u^1 du dt \right] \\
&= 2 \int_I \mathbb{E} \left[\int_I \mathbb{1}_{\{u \geq t\}} \mathbb{1}_{\{X_{t-}^1 = 1\}} \mathbb{1}_{\{X_{u-}^1 = 1\}} s(u) du \right] s(t) dt \\
&= 2 \int_I \mathbb{P}(X_{t-}^1 = 1, T_{1,0}^1 \in I, T_{1,0}^1 \geq t) s(t) dt \\
&\leq 2 \int_I \mathbb{P}(X_{t-}^1 = 1) s(t) dt = 2 \mathbb{E} \left[\int_I s(t) Y_t^1 dt \right].
\end{aligned}$$

References

- ANDERSEN, P., BORGAN, O., GILL, R. and KEIDING, N. (1993). *Statistical Models Based on Counting Processes*. Springer-Verlag, New York.
- ANTONIADIS, A. (1989). A penalty method for nonparametric estimation of the intensity function of a counting process. *Ann. Inst. Statist. Math.* **41**, 781–807.
- ANTONIADIS, A., BESBEAS, P. and SAPATINAS, T. (2001). Wavelet shrinkage for natural exponential families with cubic variance functions. *Sankhya* **63**, 309–327.
- ANTONIADIS, A. and SAPATINAS, T. (2001). Wavelet shrinkage for natural exponential families with quadratic variance functions. *Biometrika* **88**, 805–820.
- BARRON, A.R., BIRGÉ, L. and MASSART, P. (1999). Risk bounds for model selection via penalization. *Probab. Theory Relat. Fields* **113**, 301–415.
- BARRON, A.R. and COVER, T.M. (1991). Minimum complexity density estimation. *IEEE Transactions on Information Theory* **37**, 1034–1054.
- BIRGÉ, L. (1983). Approximation dans les espaces métriques et théorie de l'estimation. *Z. Wahrscheinlichkeitstheorie Verw. Geb.* **65**, 181–237.
- BIRGÉ, L. (2006). Model selection via testing : an alternative to (penalized) maximum likelihood estimators. *Ann. Inst. Henri Poincaré Probab. et Statist.* **42**, 273–325.

- BIRGÉ, L. and MASSART, P. (1998). Minimum contrast estimators on sieves: exponential bounds and rates of convergence. *Bernoulli* **4**, 329-375.
- BIRGÉ, L. and MASSART, P. (2000). An adaptive compression algorithm in Besov spaces. *Constructive Approximation* **16** 1-36.
- BIRGÉ, L. and MASSART, P. (2001). Gaussian model selection. *J. Eur. Math. Soc.* **3**, 203-268.
- BREIMAN, L., FRIEDMAN, J.H., OLSHEN, R.A. and STONE, C.J. (1984). *Classification and Regression Trees*. Wadsworth, Belmont.
- CASTELLAN, G. (1999). Modified Akaike's criterion for histogram density estimation. Technical Report 99.61. Université Paris-Sud, Orsay.
- CASTELLAN, G. (2000). Sélection d'histogrammes à l'aide d'un critère de type Akaike. *C.R.A.S.* **330**, 729-732.
- DeVORE, R.A. (1998). Nonlinear Approximation. *Acta Numerica* **7**, 51-150.
- DeVORE, R.A. and LORENTZ, G.G. (1993). *Constructive Approximation*. Springer-Verlag, Berlin.
- DeVORE, R.A. and YU, X.M. (1990). Degree of adaptive approximation. *Math. Comp.* **55**, 625-635.
- GEY, S. and NÉDÉLEC, E. (2005). Model selection for CART regression trees. *IEEE Transactions on Information Theory* **51**, 658-670.
- GRÉGOIRE, G. and NEMBÉ, J. (2000). Convergence rates for the minimum complexity estimator of counting process intensities. *J. Nonparametr. Statist.* **12**, 611-643.
- KOLACZYK, E. (1999). Wavelet shrinkage estimation of certain Poisson intensity signals using corrected threshold. *Statistica Sinica* **9**, 119-135.
- KOLACZYK, E. and NOWAK, R. (2004). Multiscale likelihood analysis and complexity penalized estimation. *Annals of Statistics* **32**, 500-527.
- LAURENT, B. and MASSART, P. (2000). Adaptive estimation of a quadratic functional by model selection. *Ann. Statist.* **28**, 1302-1338.
- LEPSKII, O.V. (1991). Asymptotically minimax adaptive estimation I: Upper bounds. Optimally adaptive estimates. *Theory Probab. Appl.* **36**, 682-697.
- MASSART, P. (2000). Some applications of concentration inequalities to Statistics. *Ann. Fac. Sciences de Toulouse* **IX**, 245-303.
- PATIL, P.N. and WOOD, A.T. (2004). A counting process intensity estimation by orthogonal wavelet methods. *Bernoulli* **10**, 1-24.
- REYNAUD-BOURET, P. (2003). Adaptive estimation of the intensity of inhomogeneous Poisson processes via concentration inequalities. *Probab. Theory Related Fields* **126**, 103-153.
- REYNAUD-BOURET, P. (2002). Penalized projection estimators of the Aalen multiplicative intensity. *School of Mathematics, Georgia Institute of Technology* Preprint 1202-002.
- STANLEY, R.P. (1999). *Enumerative Combinatorics, Vol. 2*. Cambridge University Press, Cambridge.
- van de GEER, S. (1995). Exponential inequalities for martingales, with application to maximum likelihood estimation for counting processes. *Ann. Statist.* **23**, 1779-1801.

WU, S.S. and WELLS, M.T. (2003) Nonparametric estimation of hazard functions by wavelet methods. *J. Nonparametr. Stat.* **15**, 187 – 203.

UNIVERSITÉ DE NICE SOPHIA-ANTIPOLIS, LABORATOIRE J-A DIEUDONNÉ, PARC VALROSE,
06108 NICE CEDEX 02

E-mail address: `baraud@math.unice.fr`

UNIVERSITÉ PARIS VI, LABORATOIRE DE PROBABILITÉS ET MODÈLES ALÉATOIRES, BOÎTE 188,
4 PLACE JUSSIEU, 75252 PARIS CEDEX 05

E-mail address: `lb@ccr.jussieu.fr`